

# Prediction of synthesis parameters for N, Si, Ge and Sn diamond vacancy centers using machine learning

Zhi Jiang,<sup>1</sup> Marco Peres,<sup>2,3,4</sup> Carlo Bradac,<sup>5</sup> and Gil Gonçalves<sup>1\*</sup>

1. TEMA, Department of Mechanical Engineering, University of Aveiro, Aveiro 3810-193 Portugal
2. Instituto de Engenharia de Sistemas e Computadores - Microsistemas e Nanotecnologia, Rua Alves Redol 9, 1000-029 Lisboa, Portugal
3. IPFN, Instituto Superior Técnico, University of Lisbon, Av. Rovisco Pais 1, 1049 001 Lisbon, Portugal
4. DECN, Instituto Superior Técnico, University of Lisbon, Estrada Nacional 10 (km 139.7), 2695 066 Bobadela, Portugal
5. Department of Physics and Astronomy, Trent University, 1600 West Bank Drive, Peterborough, ON, K9L 0G2, Canada

\* Corresponding author: Gil Gonçalves, [ggoncalves@ua.pt](mailto:ggoncalves@ua.pt)

**Keywords:** High Pressure High Temperature (HPHT), Microwave Plasma Chemical Vapor Deposition (MPCVD), Ion Implantation, Electron/Ion Irradiation, Diamond Color Centers, Machine Learning, Decision Tree Regression (DTR), Extreme Gradient Boosting (XGB), Shapley Values.

## Abstract

Diamond and diamond color centers have become prime hardware candidates for solid state-based technologies in quantum information and computing, optics, photonics and (bio)sensing. The synthesis of diamond materials with specific characteristics and the precise control of the hosted color centers is thus essential to meet the demands of advanced applications. Yet, challenges remain in improving the concentration, uniform distribution and quality of these centers. Here we perform a review and meta-analysis of some of the main diamond synthesis methods and their parameters for the synthesis of N-, Si-, Ge- and Sn-vacancy color-centers, including worldwide trends in fabrication techniques and processes. We extract quantitative data from over 60 experimental papers and organize it in a large database (170 data sets and 1692 entries). We then use the database to train two machine learning algorithms to make robust predictions about the fabrication of diamond materials with specific properties from careful combinations of synthesis parameters. We use traditional statistical indicators to benchmark the performance of the algorithms and show that they are powerful and resource-efficient tools for researchers and material scientists working with diamond color centers and their applications.

## **Table of contents**

1. Introduction
2. Methods and theoretical models
  - 2.1 Diamond syntheses and dopants inclusion techniques
    - 2.1.1 High pressure high temperature (HPHT)
    - 2.1.2 Chemical vapor deposition (CVD)
    - 2.1.3 Ion Irradiation
    - 2.1.4 Electron/ion irradiation
  - 2.2 Database
  - 2.3 Machine learning models and algorithms
    - 2.3.1 Decision Tree Regression (DTR)
    - 2.3.2 Extreme Gradient Boosting (XGB)
    - 2.3.2 Shapley value analysis
3. Results and discussion
  - 3.1 Task I: material synthesis
  - 3.2 Task II: Debye-Waller factor
  - 3.3 Shapley value analysis of material synthesis
4. Worldwide trends in diamond synthesis
5. Conclusions
6. Supporting Information
7. Contributions and Acknowledgements
8. Conflict of Interest
9. Data Availability
- References

## 1. Introduction

Diamond synthesis can be carried out through several fabrication processes each involving multiple variables such as pressure, temperature, chemical precursors, substrates, run time, etc. that can display complex nonlinear relations.<sup>[1]</sup> In addition, diamond can be synthesized to host extrinsic defects of foreign atoms that display unique spin-optical properties and have been successfully employed in a multitude of both fundamental and practical applications including in quantum communication<sup>[2–5]</sup> and computation,<sup>[6–8]</sup> quantum simulation,<sup>[9,10]</sup> and quantum sensing.<sup>[11–13]</sup> However, traditional orthogonal synthesis methods often struggle to effectively capture and single out the optimal experimental parameters that lead to the synthesis of diamond materials and their dopant concentrations with desired characteristics. This limitation opens the opportunity for alternative approaches such as machine learning (ML), which shows incredible promise in contexts concerning the effective extraction of patterns from complex data and the capture of difficult-to-model physical mechanisms akin to the ones involved in material synthesis.<sup>[14,15]</sup> Models based on ML algorithms can, for instance, quickly identify the key parameters of fabrication processes<sup>[16–18]</sup>, and can be used to tailor and design specific material properties,<sup>[19]</sup> heterostructures,<sup>[20,21]</sup> and devices.<sup>[22–25]</sup>

Currently, in material science the focus is on the synthesis of diamond hosting extrinsic dopants, specifically on controlling the inclusion of diamond vacancy color centers (DVCCs). These are atom(s)-vacancy(ies) complexes replacing carbon atoms in the diamond crystalline matrix: the most prominent examples are the nitrogen-vacancy (NV) center and the group IV diamond color centers (SiV, GeV, SnV and PbV).<sup>[26,27]</sup> Notably, achieving precise control of these vacancy color centers, whilst fabricating high-quality diamond films and nanomaterials, remains a significant challenge under current technological capabilities.<sup>[28–32]</sup>

In this study, we propose the use of ML to address this challenge. Our approach is based on two ML methods: Decision Tree Regression (DTR) and Extreme Gradient Boosting (XGB).<sup>[33–35]</sup> Briefly, we perform an extensive meta-analysis of the existing literature (60 studies) and create a large database (170 data sets and 1692 entries) containing the synthesis parameters and color centers yield for each of the identified diamond fabrication techniques. We then use this database to train-and-test the DTR and XGB algorithms and perform Shapley value analyses to identify the optimal combinations of synthesis parameters that produce diamond materials with desired characteristics and high-quality color centers. Specifically, we train and optimize the DTR and XGB algorithms with the constructed database to predict the *target value* of certain variables of interest from the values of all the other *input features* in the database. We demonstrate the effectiveness of this approach for each synthesis technique by performing two *Tasks*. *Task I* (§ 3.1) is specific to each synthesis method: the target variable and input features are chosen amongst the relevant synthesis parameters of the method. For example, we use the DTR/XGB algorithms to predict the size (target variable) of diamond materials synthesized by

high pressure high temperature (HPHT) from the value of the chamber temperature, chamber pressure and run time (input features) of the process. *Task II* (§ 3.2) is more general and comparative: for all techniques, we select the Debye-Waller factor (DWF) of the DVCCs as the target variable, while the input features are chosen amongst the relevant synthesis parameters of each method. This allows us to directly compare all synthesis techniques. For *Task II*, we use the DWF as a proxy for the quality of the DVCCs. This is motivated by a series of observations: *i)* the need to base our analysis on a *measurable* physical observable that is *ii)* *common* to all synthesis techniques, and *iii)* that is *physically meaningful* (many of the technological quantum realization involving DVCCs heavily rely on their optical and spectral properties, for which the DWF can be a useful indicator).

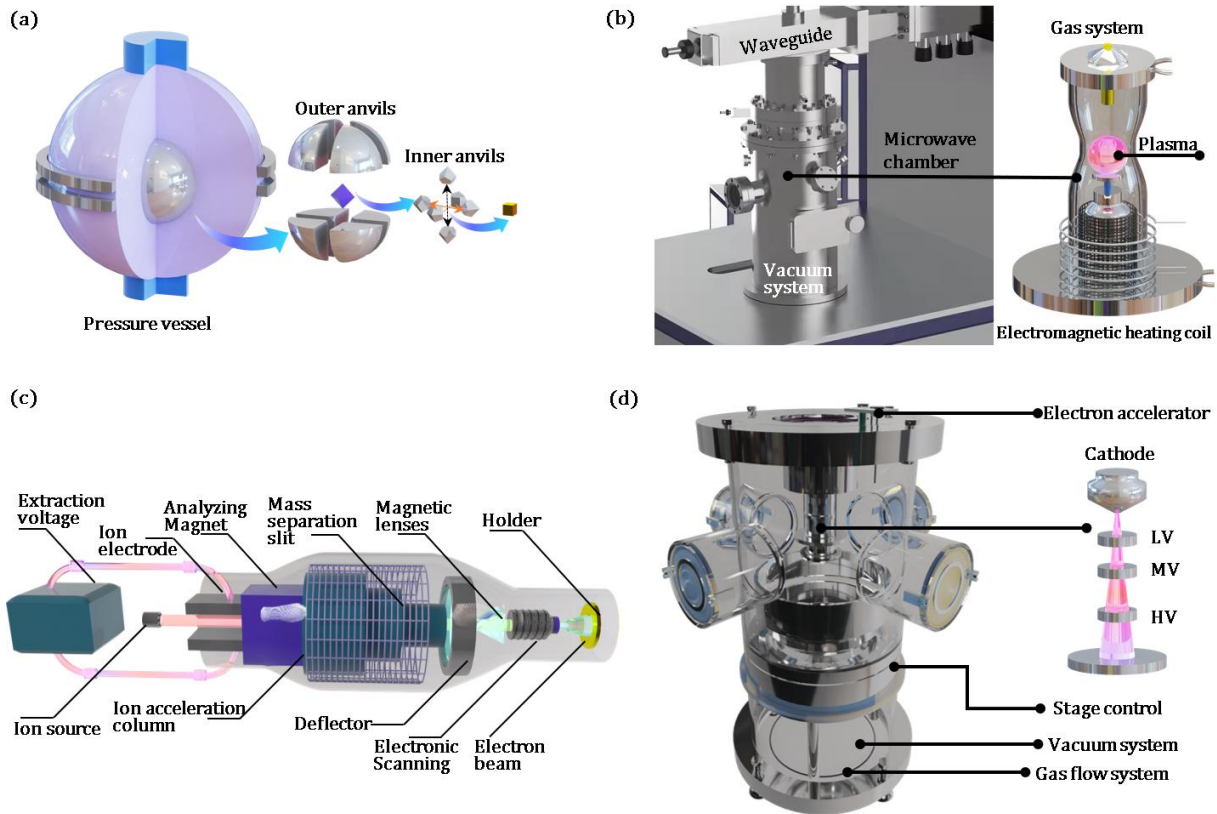
The goal of this study is thus to provide a deterministic and practical toolset for diamond synthesis and design that focuses on the yield of materials with optimal optical properties and performance.

## 2. Methods and theoretical models

Our analysis focuses on four different processes for synthesizing diamond materials and controlling the inclusion of extrinsic defects, particularly those forming DVCCs. We concentrate on high pressure high temperature (HPHT) and chemical vapor deposition (CVD) for material synthesis, and on ion implantation and electron/ion irradiation for the control and inclusion of target dopants. These techniques are briefly reviewed in §§ 2.1.1–2.1.4 to establish context, help identify the critical parameters of each technique, and detail the implementation of *Tasks I* and *II* for each individual method. The construction of the database is briefly discussed in § 2.2 (more details can be found in the Supplementary Information, SI, § S1). The methods and ML algorithms employed to extract the values of said parameters for the optimal synthesis of diamond and color centers are discussed in § 2.3.

### 2.1 Diamond synthesis and dopants inclusion techniques

The first step of our analysis consists in selecting a set of techniques of interest for diamond synthesis and the inclusion of dopants in the fabricated materials. These methods are summarized in Figure 1 and are reviewed in the following §§ 2.1.1–2.1.4, with the main intent of highlighting, for each technique, the experimental variables that affect the yield of material and of DVCC inclusion. This is important because we use these variables in the database to train the proposed machine learning DTR and XGB algorithms and make predictions of the identified target values from the input features. To compile the database, we reviewed over 60 published papers, gathering a total of 170 datasets, summarized and organized in Tables ST2–ST5 (see SI, § S2).



**Figure 1.** **A)** HPHT setup. Carbon-based precursor materials are subjected to high pressures and high temperatures through multiple sets of nested anvils controlled by hydraulic, electrical and heating/cooling systems. **B)** MPCVD setup. Diamond is grown atom-by-atom in a vacuum chamber from a hydrocarbon gas mixture. The setup includes a deposition chamber, a gas delivery system, a microwave source, a vacuum system and a cooling system. **C)** Ion implantation setup. Ions are accelerated and implanted on a target substrate. The setup includes an ion source and accelerator, a target chamber, a vacuum system, beam deflection and scanning systems, and a dose and energy control system. **D)** Electron/ion irradiation setup. A focused beam of electrons is used to irradiate the target diamond material and create vacancies to form DVCCs. The setup includes an electrons gun with voltage acceleration anodes, as well as focusing and deflection systems and a vacuum chamber for the target.

### 2.1.1. High Pressure High Temperature (HPHT)

The synthesis of artificial diamond by high pressure high temperature (HPHT) finds its origin in the mid-1950s.<sup>[36]</sup> In this process, which mimics the geological formation of natural diamond, carbon-based precursor materials are subjected to high temperature (above 1500 °C) and high pressure (above 5 GPa) to induce the phase transition of graphite-like carbon to diamond. This technique is usually associated with high process and equipment costs, high energy consumption, and limitations in large-scale production,<sup>[37]</sup> but it reliably yields high-quality diamond material and large-grained, single-crystalline micro/nanoparticles (M/NPs) capable of hosting color centers. Notably, and interestingly for this work, the HPHT process is attractive as it allows for the inclusion of DVCCs directly during synthesis. For instance, the employment

of inorganic or organic compounds high in nitrogen content, such as  $\text{NaN}_3$ ,  $\text{C}_3\text{H}_6\text{N}_6$ ,  $\text{Ba}(\text{N}_3)_2$  and  $\text{Fe}_3\text{N}$ , can lead to the relatively high incorporation of nitrogen-vacancy (NV) centers directly during material fabrication.<sup>[38–44]</sup>

The characteristics of synthesized HPHT diamond depend on its fabrication parameters. For our analysis, we identified five key variables that prominently determine the yield and properties of HPHT diamond and its color centers: chamber pressure and temperature, run time, diameter of the synthesized diamonds and Debye-Waller factor (DWF) of the hosted DVCCs. With relevance to the *Tasks* introduced in § 1, we have the following for our analysis (see SI, § S1, Table ST1). For *Task I*, the target variable is the average diameter of the synthesized diamond NPs, and the input features are the chamber pressure and temperature, and the run time. For *Task II*, the target variable is the DWF of the DVCCs, and the input features are the chamber pressure, chamber temperature, run time, and average diameter of the synthesized diamond NPs.

### 2.1.2. Chemical vapor deposition (CVD)

Chemical vapor deposition (CVD) and microwave plasma chemical vapor deposition (MPCVD) are diamond synthesis techniques that consist in the atom-by-atom growth of diamond on a heated substrate from a chemically reacting hydrocarbon gas mixture. Usually, the growth process is carried out inside a vacuum chamber where temperature, pressure and, in the case of MPCVD, microwave (MW) frequency and power are precisely controlled. Research on the growth of artificial diamonds from the gas phase started in the early 1950s,<sup>[45]</sup> with other major milestones in the development of the technique following in the late 1960s,<sup>[46]</sup> 1970s<sup>[47]</sup> and 1980s.<sup>[48–51]</sup> Diamond synthesis via CVD is quite an attractive method as it can yield materials with a certain flexibility in terms of size and properties, whilst also offering practical advantages. These include the lack of need for extremely high pressures (which can be limiting in HPHT techniques, see § 2.1.1 above) and, in the case of MPCVD, the fact that the material synthesis process takes place at relatively high growth rates, efficient carbon source utilization, lower deposition temperatures and low contamination.<sup>[52–54]</sup> Similarly to the HPHT process, CVD growth also allows for the inclusion of color centers in the diamond material directly during synthesis by injection in the chamber of gas mixtures containing specific dopants. It should however be remarked that the incorporation of dopants during growth can be quite inefficient,<sup>[55–58]</sup> and it requires careful control over experimental variables such as chemical composition of the reaction atmosphere (e.g. via the use of inert gases such as  $\text{N}_2$  and Ar, and the reduction of  $\text{H}_2$ ), pressure, substrate temperature, and power of the MW radiation.<sup>[59]</sup>

The synthesis of diamond by CVD and MPCVD thus involves a complex and large set of inter-dependent physical mechanisms and chemical reactions. To build our database and train our algorithms to make quantitative predictions about quality, rate, morphology, doping, and DVCC inclusion of the synthesized diamond, we use the following parameters: chamber pressure,

substrate temperature, MW plasma power, and DWF of the hosted DVCCs. With relevance to the *Tasks* introduced in § 1, we have the following for our analysis (see SI, § S1, Table ST1). For *Task I*, the target variable is the MW plasma power, and the input features are the chamber pressure and substrate temperature. For *Task II*, the target variable is the DWF of the DVCCs, and the input features are the chamber pressure, substrate temperature, and MW power.

### 2.1.3. Ion implantation

Ion implantation is a process for the inclusion of foreign atoms in the diamond crystalline matrix after synthesis.<sup>[60]</sup> In this process, high-energy ions of specific elements are accelerated towards the target diamond and are embedded into it.<sup>[61–63]</sup> Ion implantation is advantageous for DVCC inclusion for it has a relatively high spatial resolution, both transversal and longitudinal, and allows for the fabrication of isolated color centers. The deterministic implantation is desirable for quantum applications as these often rely on the ability to isolate individual single-photon sources.<sup>[64–67]</sup> Further control over the creation of DVCCs can be achieved through low-energy carbon ion implantation, which can induce the formation of atom-vacancy complexes in the shallow surface layer of diamond, without introducing additional impurities.<sup>[68–71]</sup>

Whilst conceptually straightforward, ion implantation is technically quite challenging. For instance, the implantation parameters for each of the different dopant species forming DVCCs (i.e. N, Si, Ge, Sn, Pb) must be controlled with high precision as they are specific to and different for each element, while also depending on the concentration of impurities created during synthesis.<sup>[38,72–77]</sup> Additionally, the stochastic nature of the implantation and the straggle of the ions require careful control over their implantation kinetic energies and fluence (i.e. the number of striking atoms per unit area) to minimize lateral displacement while maximizing DVCCs yield.<sup>[76,78]</sup> Ion implantation usually also involves at least one annealing step (at temperatures  $\geq 850$  °C). The reason is twofold: to repair the damage to the diamond  $sp^3$  bonds caused by the implanted ions,<sup>[66,71,74,75]</sup> and to activate the mobilization of vacancies and thus promote the formation of DVCCs in the diamond material.<sup>[79,80]</sup>

Obtaining diamond materials with desired characteristics and properties through ion implantation is markedly non-trivial, for the combinatorics of the values for the synthesis parameters is considerable. To build our database and train our algorithms, we use the following parameters for ion irradiation: implantation energy, fluence, annealing temperature, annealing time and DWF of the hosted DVCCs. With relevance to the *Tasks* introduced in § 1, we have the following for our analysis (see SI, § S1, Table ST1). For *Task I*, the target variable is the annealing time, and the input features are the implantation energy, fluence, and annealing temperature. For *Task II*, the target variable is the DWF of the DVCCs, and the input features are the implantation energy, fluence, annealing temperature, and annealing time.

#### 2.1.4. Electron/ion irradiation

Electron/ion irradiation is a common approach to induce the formation of DVCCs in diamond. A focused beam of electrons (or ions) is used to irradiate the target diamond material and create vacancies in the crystalline structure. Subsequent annealing at temperatures  $\geq 850$  °C promotes the migration of the vacancies, which become trapped by the foreign atoms already present in the diamond lattice, forming stable DVCCs.<sup>[81–83]</sup> The annealing step also helps repair the damage caused by the irradiation process. This is particularly relevant for this study as the optical and spectroscopic properties of the DVCCs are affected by strain, which can alter the photon purity,<sup>[84,85]</sup> lifetime<sup>[86]</sup> and coherence<sup>[87]</sup> of the DVCCs. Ad-hoc annealing conditions and advanced annealing strategies (e.g. laser and pulsed annealing<sup>[88,89]</sup>) can thus strongly affect the number and quality of the DVCCs and are taken into consideration in our models.

To build our database and train our algorithms, we use the following parameters for electron/ion irradiation: electron/ion irradiation energy, fluence, annealing temperature, annealing time, and DWF of the hosted DVCCs. With relevance to the *Tasks* introduced in § 1, we have the following for our analysis (see SI, § S1, Table ST1). For *Task I*, the target variable is the annealing time, and the input features are the electron/ion irradiation energy, fluence, and annealing temperature. For *Task II*, the target variable is the DWF of the DVCCs, and the input features are the irradiation energy, fluence, annealing temperature, and annealing time.

## 2.2 Database

The next step in our analysis is the construction of the database. We use the database to train and test the DTR and XGB algorithms and, for each synthesis technique, make predictions about the selected target variables from the input features. The database is compiled from 60 experimental studies and has a total of 1692 entries, organized in 170 datasets and 4 sub-databases covering each one of the 4 synthesis methods; this information is displayed and organized, for clarity, in Tables ST2–ST5 (see SI, § S2). Each table is arranged by color center: NV, SiV, GeV and SnV, and shows the values and specifications for the critical synthesis parameters discussed in secs. 2.1.1–2.1.4, with regard to *Tasks I* and *II*.

## 2.3 Machine learning models and algorithms

In this work, we use two machine learning algorithms to make predictions about target values from input features: Decision Tree Regressor (DTR) and Extreme Gradient Boosting (XGB). We also implement Shapley value analysis to quantify the contributions of the various input features to the target variables. These methods are briefly reviewed below from a theoretical standpoint; their application to the collected data is discussed in §§ 3.1–3.3.



### 2.3.1. Decision Tree Regression (DTR)

Decision Tree Regression (DTR) is a supervised learning algorithm that models the relation between a set of input features and a continuous target variable.<sup>[90]</sup> The model constructs a decision tree by recursively partitioning the data into subsets based on the values assumed by the input features, aiming to minimize the prediction error within each subset. In our case the data is organized as follows. Each one of the 4 sub-databases has a total number of observations,  $n$ . These are the rows in the sub-database (SI, § S2, Tables ST2–ST5) constructed from the meta-analysis of the literature. Each diamond synthesis method has a number of input features,  $d$  (summarized for each technique and for each *Task* in Table 1), and a target variable,  $y$ . The target variable is the observable whose value ( $\in \mathbb{R}$ ) we are interested in predicting and is specified for each technique and for each *Task* in the SI, § S1, Table ST1.

Effectively, the data has form  $(\mathbf{X}, \mathbf{y})$  where  $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]^T \in \mathbb{R}^{n \times d}$  is the feature matrix and  $\mathbf{y} = [y_1, y_2, \dots, y_n]^T \in \mathbb{R}^n$  is the target vector, such that each  $i$ -th observation  $\mathbf{x}_i = [x_{i1}, x_{i2}, \dots, x_{id}] \in \mathbb{R}^d$  is a row of input features with corresponding output  $y_i \in \mathbb{R}$ . For each synthesis method, the algorithm starts from the entire dataset considering all the features and proceeds in steps such that, at each steps, it finds the best feature  $j$  and corresponding threshold  $s$  with respect to which the data should be split in left,  $L: \{\mathbf{x}_i \mid x_{ij} \leq s\}$ , and right,  $R: \{\mathbf{x}_i \mid x_{ij} > s\}$ . The best split is defined with respect to the criterion of minimizing the sum of squared errors (SSE) in each  $L$  and  $R$  partition:

$$\text{SSE}_{\text{split}} = \sum_{i \in L} (y_i - \bar{y}_L)^2 + \sum_{i \in R} (y_i - \bar{y}_R)^2 \quad (1)$$

Where  $\bar{y}_L$  and  $\bar{y}_R$  are the mean target values of the  $L$  and  $R$  partitions, respectively. The algorithm proceeds by splitting the data in each subset recursively until a stopping criterion is met (e.g., minimum number of samples, or no gain in error reduction). In the final tree, each leaf node predicts a constant value, which is the mean of all  $y_i$  values in that region.

Once the regressor tree is constructed, if an unknown input  $\mathbf{x}_u = [x_{u1}, x_{u2}, \dots, x_{ud}] \in \mathbb{R}^d$  is passed to it, the tree routes it through a set of decisions until it reaches a leaf node of output:

$$\hat{y}_u = \frac{1}{|S|} \sum_{i \in S} y_i \quad (2)$$

Where  $S$  is the set of training samples in that leaf region.

### 2.3.2. Extreme Gradient Boosting (XGB)

Extreme Gradient Boosting (XGB) is a supervised learning algorithm that combines the predictions of many decision trees to produce a more accurate model than a single tree (see DTR above).<sup>[91]</sup> It uses an ensemble technique called gradient boosting, which builds trees

sequentially, with each tree learning to correct the residual errors of the previous ones. As for the case of DTR, the input features and a target variable are identified for each set of data of interest (in our case, for each synthesis technique and each *Task*).

The data is formatted in the same way as per the DTR case (see § 2.3.1), in the form  $(\mathbf{X}, \mathbf{y})$ , where  $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]^T \in \mathbb{R}^{n \times d}$  is the feature matrix and  $\mathbf{y} = [y_1, y_2, \dots, y_n]^T \in \mathbb{R}^n$  is the target vector, such that each  $i$ -th observation  $\mathbf{x}_i = [x_{i1}, x_{i2}, \dots, x_{id}] \in \mathbb{R}^d$  is a row of input features with corresponding output  $y_i \in \mathbb{R}$ . As per the DTR case,  $n$  is the number of observations,  $d$  is the number of features and  $y$  is the target variable. The XGB algorithm builds the model as a sum of many regression trees:

$$\hat{y}_i = \sum_{k=1}^K f_k(\mathbf{x}_i), \quad f_k \in \mathcal{F} \quad (3)$$

Where  $\hat{y}_i$  is the prediction for the  $i$ -th sample,  $f_k$  is the  $k$ -th tree function (decision tree),  $\mathcal{F}$  is the space of regression trees and  $K$  is the number of boosting rounds (trees). Each new tree is added to minimize a loss function, which measures the error between predicted and actual values. Specifically, the algorithm aims at minimizing the object:

$$\mathcal{L}^{(t)} = \sum_{i=1}^n \ell(y_i, \hat{y}_i^{(t-1)} + f_t(\mathbf{x}_i)) + \Omega(f_t) \quad (4)$$

Where  $\ell$  is the loss function (e.g., the squared error  $\ell(y, \hat{y}) = (y - \hat{y})^2$ ), the term  $\hat{y}_i^{(t-1)}$  is the current prediction,  $f_t$  is the new tree and  $\Omega(f_t)$  is the regularization term, i.e. a function that represents a penalty for the complexity of a tree and discourages overly complex trees (briefly, trees that have more leaves and are deeper can model data very closely, but run into the risk of overfitting the data, which this regularization term helps avoid). To fit the new tree, XGB uses Taylor expansion of the loss function (gradient boosting). Unlike DTR, which builds a single deep tree, XGB builds shallower trees and adds them up. Each tree is trained to focus on the residuals (errors) from the previous model. After all  $K$  trees are trained, any unknown input  $\mathbf{x}_u = [x_{u1}, x_{u2}, \dots, x_{ud}] \in \mathbb{R}^d$  passed to the algorithm produces the target variable prediction:

$$\hat{y}_u = \sum_{k=1}^K f_k(\mathbf{x}_u) \quad (5)$$

Which is usually a more robust, less overfitted model than using one tree (DTR).

To train and test our DTR and XGB algorithms we implemented a common 70-30 split of the data. That is, 70% of the database entries are randomly used for training the algorithms, and the remaining 30% are used to test the data. The errors are calculated on the 30% test data and are given as the differences between the true experimental values in the database and the values produced by the DTR/XGB predictors (see § 3 for details). A schematic implementation of the XGB algorithm to the synthesis of diamond material via HPHT, used as an example, is shown in the SI, § S1, Fig. SF1. The list of hyperparameters used in our models and the range of their

values can instead be found in the SI, § S3, Table ST6.

We note that in our analysis, we excluded certain *categorical* parameters for the optimization of the models, such as the carbon source for HPHT and CVD/MPCVD, the dopant source and substrate for CVD/MPCVD, the ion source for Ion implantation and Electron/ion irradiation, etc. While in principle these affect the synthesized materials, we excluded them for practical, defensible reasons. These variables are difficult to include in our models as they would require qualitative, categorical encoding, i.e., converting non-numeric, qualitative data into numeric form. Additionally—and more importantly—a post-facto analysis of our results shows that already with the input features listed in Table ST1 (SI, § S1), the optimized DTR/XGB models have good predictive power (coefficient of determination  $R^2 \sim 1$ , and relatively small mean squared error (MSE) and mean absolute error (MAE), see §§ 3.1, 3.2), without the need to include these categorical features—thus justifying their exclusion from the model.

### 2.3.3. Shapley value analysis

In addition to using DTR and XGB to make predictions, we use Shapley value analysis to establish the relevance of each parameter in determining the characteristics of the synthesized diamond materials.<sup>[92]</sup> Shapley value analysis is a method developed from cooperative game theory that helps explain the contribution of each feature to a target value prediction made by a (machine learning) model, through fair, additive, and consistent feature attributions.

As per our discussion in §§ 2.3.1 and 2.3.2, the data is in the form  $(\mathbf{X}, \mathbf{y})$ , where  $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]^T \in \mathbb{R}^{n \times d}$  is the feature matrix and  $\mathbf{y} = [y_1, y_2, \dots, y_n]^T \in \mathbb{R}^n$  is the target vector, such that each  $i$ -th observation  $\mathbf{x}_i = [x_{i1}, x_{i2}, \dots, x_{id}] \in \mathbb{R}^d$  is a row of input features with corresponding output  $y_i \in \mathbb{R}$ . Again,  $n$  is the number of observations,  $d$  is the number of features and  $y$  is the target variable. For each individual observation,  $\mathbf{x}_i = [x_{i1}, x_{i2}, \dots, x_{id}]$ , a trained model  $f$  gives a prediction:

$$\hat{y}_i = f(\mathbf{x}_i) \quad (6)$$

Where in this equation,  $f$  is either the DTR or XGB regressor we built from the available datasets and whose explicit expressions are given in equations (2) and (5), respectively. Shapley value analysis decomposes this prediction into a sum of feature contributions:

$$f(\mathbf{x}_i) = \phi_0 + \sum_{j=1}^d \phi_j^{(i)} \quad (7)$$

Where  $\phi_0$  is the baseline prediction (often the mean prediction over the training set), and  $\phi_j^{(i)}$  is the  $j$ -th Shapley value, i.e., the contribution of the  $j$ -th feature to the prediction for the  $i$ -th *instance*. The Shapley value calculates the average marginal contribution of a feature across all possible subsets of other features for each  $i$ -th instance. Effectively, it is the average change in

the model output when feature  $j$  is added to all possible subsets of other features. Mathematically, for feature  $j$ , the Shapley value for  $\mathbf{x}_i$  is:

$$\phi_j^{(i)} = \sum_{S \subseteq \mathcal{F} \setminus \{j\}} \frac{|S|!(d-|S|-1)!}{d!} [f_{S \cup \{j\}}(\mathbf{x}_i) - f_S(\mathbf{x}_i)] \quad (8)$$

Where  $\mathcal{F} = \{1, 2, \dots, d\}$  is the full set of features,  $S$  is a subset of features excluding the  $j$ -th, and  $f_S(\mathbf{x}_i)$  is the model's prediction when only features in  $S$  are known; the weighting factor in front ensures fairness over all permutations.

The Shapley value analysis thus allows us, for DTR and XGB, to add numerical, per-feature attribution to predictions, effectively helping us quantify how much each feature contributed for each individual sample. This can be seen explicitly in equation (7) where the prediction  $f(\mathbf{x}_i)$  is decomposed in a sum of terms,  $\phi_0$  and  $\phi_j^{(i)}$ , where each term  $\phi_j^{(i)}$  has a numerical value (calculated using eq. (8)) that can be compared to that of any other in the sum. The quantitative Shapley value analysis applied in this study is discussed in detail in § 3.3.

### 3. Results and discussion

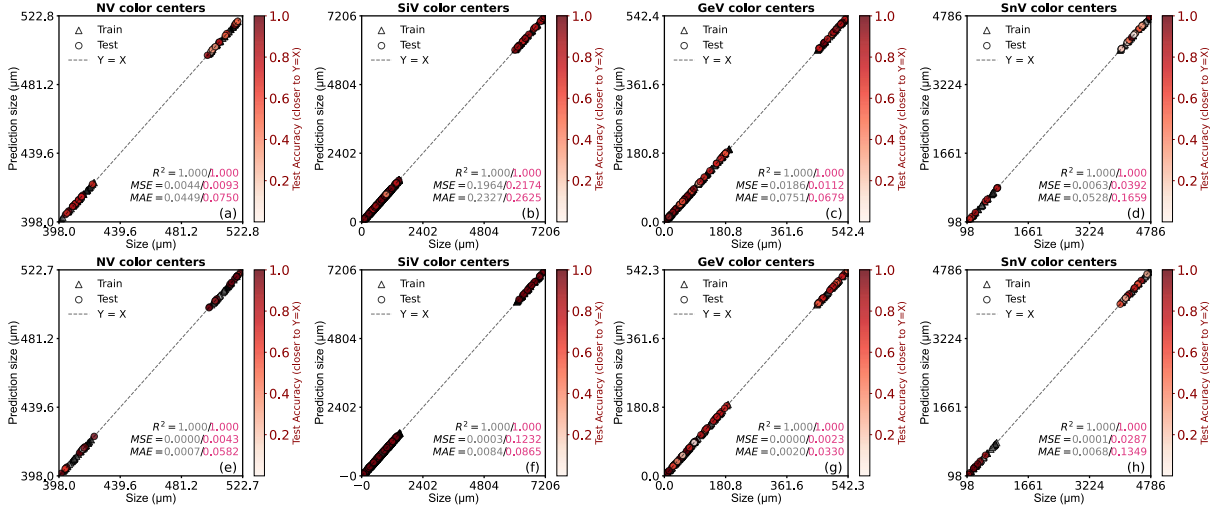
In this section, we present and discuss the effectiveness of our DTR and XGB algorithms in predicting target variables from input features. The analysis is organized with respect to the *Tasks I* and *II* we identified for each technique in §§ 2.1.1–2.1.4 and Table ST1 (SI, § S1).

#### 3.1. Task I: material synthesis

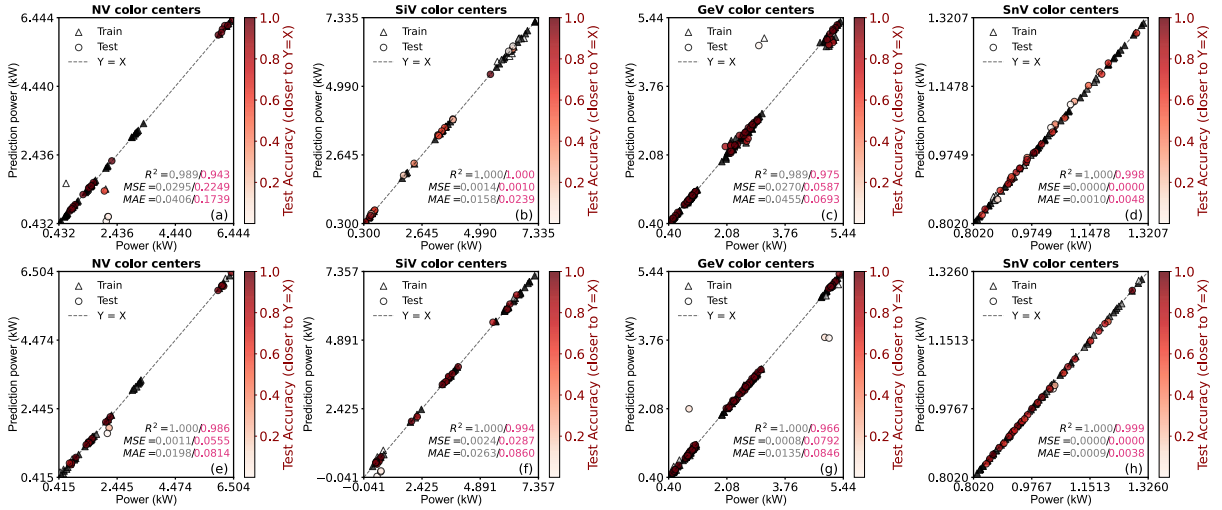
To measure the effectiveness of the DTR and XGB predictors quantitatively, we use traditional statistical estimators. Training and testing of the models follow a typical 70-30 split with the predictors been trained on 70% of the database entries, at random, and tested on the remaining 30%. Accuracy and precision of the predictors are determined from the errors, which are calculated based on the differences between the true and predicted values. The true values of the target variables are those reported in the experimental data of the database. The predicted values are instead the values of the target variables as estimated by the DTR and XGB predictors. We quantified the performance of the predictors by measuring their coefficient of determination ( $R^2$ ), mean squared error (MSE) and mean absolute error (MAE).

Figures 2–5 show the main results for *Task I* of each technique organized by diamond color center (columns) and by predictor (rows). The  $R^2$ , MSE and MAE values for each technique are displayed directly in the figure subpanels both for the test (circles) and training sets (triangles); the values of the latter are merely used as control. The values of  $R^2$ , MSE and MAE for the DTR/XGB predictors for each synthesis method are also listed in the SI, § S4, Table ST7. In Figs. 2–5, the  $x$ -axis is the true value of the target variable and the  $y$ -axis the predicted one.

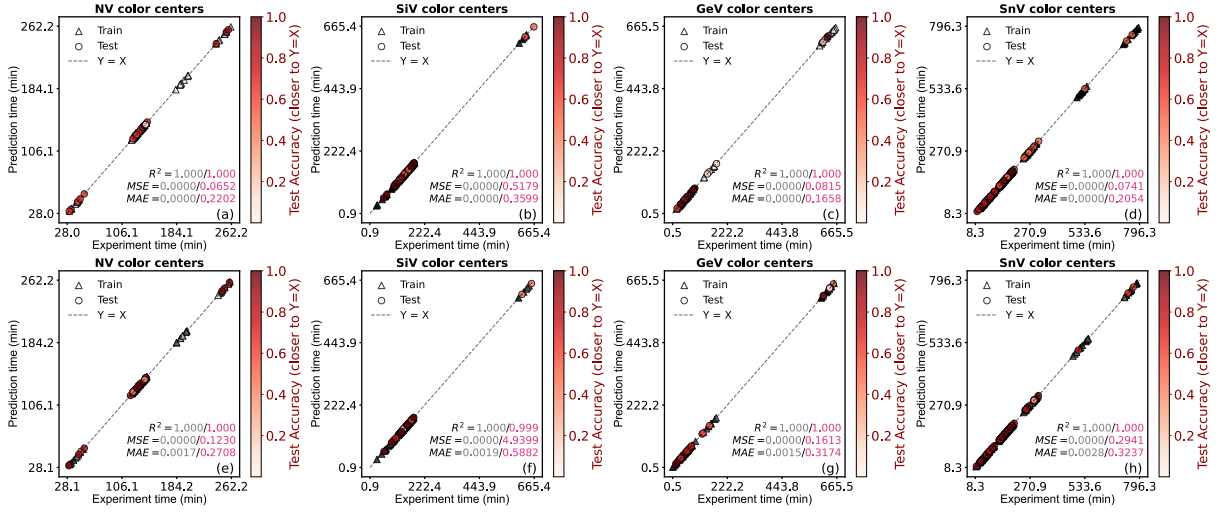
The diagonal is thus the *line of perfect accuracy*: a point on this line indicates that the DTR/XGB predictors estimated a target value from the input features coinciding exactly with the true experimental value. Points above and below the line correspond to the model over- and under-predicting the true value, respectively. The distance of the points from this line is thus a measure of the predictors' accuracy (for clarity, in Figs. 2–5, we also use color-coding: the higher the accuracy, the darker the color), while the MSE/MSA values are indicative of their precision.



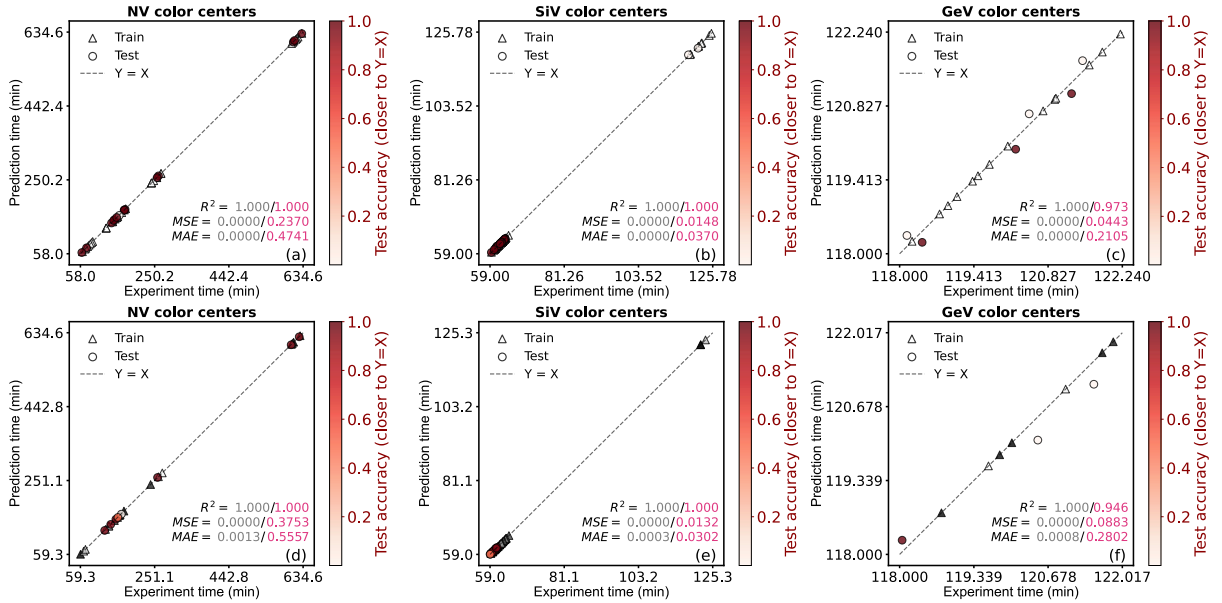
**Figure 2.** HPHT, *Task 1* summary. DTR (top row) and XGB (bottom row) prediction of the average HPHT diamond diameter (target variable) from given values for chamber pressure, chamber temperature, and run time (input features), organized by color centers (columns). The values for  $R^2$ , MSE and MAE for the training (triangles) and test (circles) sets are listed in the bottom-right corner of each subpanel.



**Figure 3.** CVD/MPCVD, *Task 1* summary. DTR (top row) and XGB (bottom row) predictions of the MW plasma power (target variable) from given values for chamber pressure and substrate temperature (input features), organized by color centers (columns). The values for  $R^2$ , MSE and MAE for the training (triangles) and test (circles) sets are listed in the bottom-right corner of each subpanel.



**Figure 4.** Ion implantation, *Task I* summary. DTR (top row) and XGB (bottom row) prediction of the annealing time (target variable) from given values for implantation energy, fluence and annealing temperature (input features), organized by color centers (columns). The values for  $R^2$ , MSE and MAE for the training (triangles) and test (circles) sets are listed in the bottom-right corner of each subpanel.



**Figure 5.** Electron/ion irradiation, *Task I* summary. DTR (top row) and XGB (bottom row) prediction of the annealing time (target value) from given values for irradiation energy, fluence and annealing temperature (input features), organized by color centers (columns). The values for  $R^2$ , MSE and MAE for the training (triangles) and test (circles) sets are listed in the bottom-right corner of each subpanel.

The numerical results and the plots of our analysis for the DTR and XGB predictors warrant a series of observations.

Overall, the values of the statistical estimators, specifically  $R^2 \sim 1$ , and the relatively small values for MSE and MAE (both for the training and test sets) indicate that the DTR and XGB models are good generalizations. Interestingly, when comparing directly the performance of the DTR

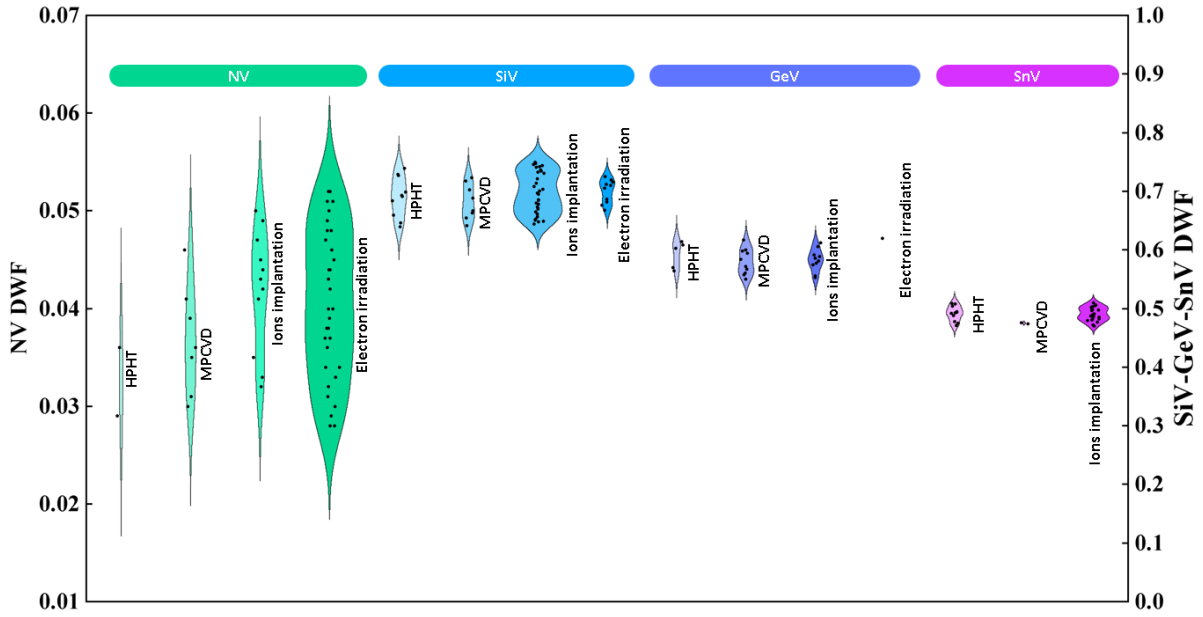
and XGB algorithms, we found that DTR is a better predictor than XGB as, in most cases, it has values for  $R^2$  closer to 1, and smaller values (roughly by a factor  $\times 5$ – $10$ ) for MSE and MAE. This is likely due to several factors.<sup>[93,94]</sup> In terms of number of input features, our data is considered of low dimensionality (only 3–4 features) compared to scenarios of medium ( $10^1$ – $10^2$  features) or high ( $10^2$ – $10^3$  features) dimensionality. In low-dimensionality problems, DTR models can learn patterns quickly and the ensemble averaging effects characteristic of XGB can become negligible. These effects are further emphasized by the fact that in our analysis the features/data-size ratio is relatively small (for us this ratio is in the low range 0.004–0.03, as we have 3–4 input features and 128–680 rows of data, depending on the database used; see SI, § S2, Tables ST2–ST5). Finally, a foreshadowing glance at the Shapley value analysis (see §§ 3.3 and Figs. 11–14 for details) indicates that in all four cases (HPHT, CVD/MPCVD, Ion implantation and Electron/ion irradiation), one feature tends to outweigh the others in explaining the variance of the target variable. When this occurs, DTR can efficiently exploit this dominance by making the first split of the tree on that feature and focusing most of the optimization on it, while using the other features for fine tuning. Conversely, XGB builds one tree at a time and such that each subsequent tree learns to fix the residuals of the previous ones. When one feature is dominant, XGB may dilute its effect by spreading influence across more, smaller steps (especially if the learning rate is low or the depth of the tree is small). This is usually desirable for generalization, but it may slow down convergence or—relevantly here—underuse dominant feature(s) signals.

Analysis of the plots in Figures 2–5, also reveals that for each technique there are certain areas of the parameter space that are unexplored or, in other words, there are no data points corresponding to certain combinations of parameters (see, e.g., the ranges 410–480  $\mu\text{m}$  in Figs. 2a and 2e, 1,800–5,000  $\mu\text{m}$  in Figs. 2b and 2f, and 150–350  $\mu\text{m}$  in Figs. 2c and 2g). This is due to a series of factors including the specific goals of the authors in the original studies, technical limitations, or the fact that, while relatively large ( $\sim 60$  papers), the number of studies we reviewed is just a subset of the existing literature. This is important because, in general, the reliability of trained models should not be trusted beyond the ranges spanned by the training datasets. However, our plots show reliable results, consistently, for different non-contiguous ranges of parameters. This is favorable, for it indicates that the interaction between the various input features and target variables is predictable and mostly devoid of interactions that our models are unable to capture. We thus expect our models to be robust predictors even with parameters in the empty areas of the parameter space—at least within the lowest and highest extremes of the explored ranges.

Finally, we note that for the cases of Ion implantation and Electron/ion irradiation (Figs. 4 and 5), we opted to preprocess the data by grouping the entries in the database by type of ion source as we observed that defecting to do so led to relatively larger test prediction errors.

### 3.2. Task II: Debye-Waller factor

We now want to generalize our approach and apply the DTR and XGB algorithms simultaneously to all synthesis methods, with the intent of drawing direct comparisons. The goal of *Task II* of our analysis is to train the DTR/XGB models to predict the value of a target variable common to all methods from the input features of each technique. The target variable of choice is the Debye-Waller factor (DWF) of the DVCCs, which we chose based on arguments discussed in § 1. The DWF quantifies the fraction of photons emitted by an emitter into the zero-phonon line (ZPL) versus those emitted into phonon sidebands (PSB), in an excited vibrational state. It is defined as the integrated luminescence intensity of the ZPL,  $I_{\text{ZPL}}$ , divided by the total integrated luminescence intensity of the color center  $I_{\text{tot}}$ <sup>[95,96]</sup> A high value of the DWF is usually desirable for optical and quantum applications based on diamond color centers as it improves photon indistinguishability, photon extraction, and the coupling efficiency to optical cavities and photonics structures.<sup>[27,97,98]</sup> We thus use it here as a figure of merit to measure the quality of the diamond color centers. Ultimately, our intent is to combine the results from *Tasks I* and *II* for the DTR/XGB predictors and identify the synthesis techniques and their parameters that yield diamond with desired material properties and DVCCs quality.

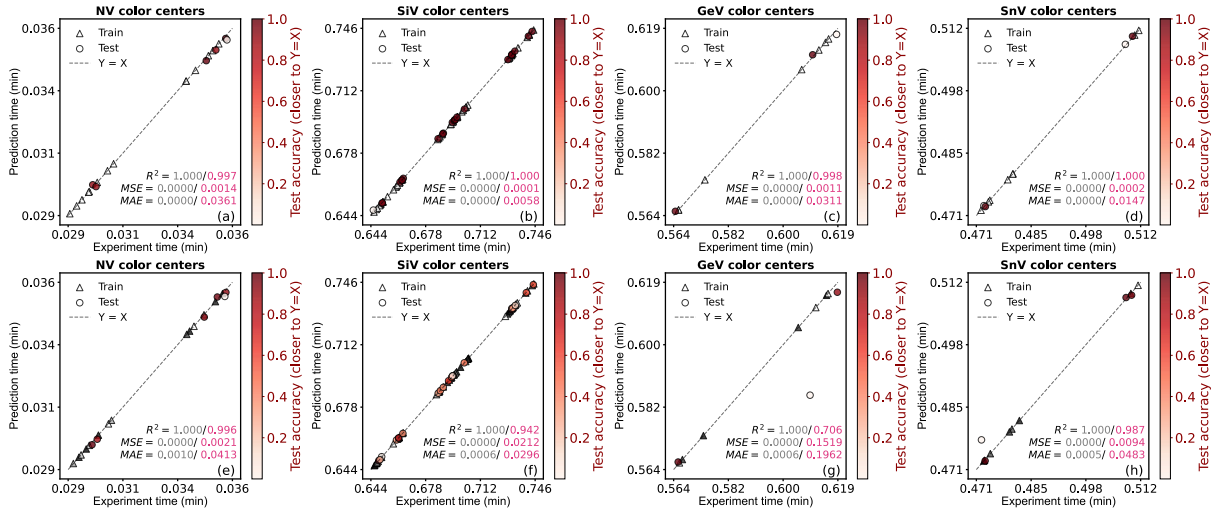


**Figure 6.** Debye-Waller factor (DWF) for NV, SiV, GeV and SnV, organized by synthesis method. The DWF values were calculated from the photoluminescence spectra available in the published experimental data. Note that the scales are different for the NV center (left y-axis) and for the SiV, GeV and SnV centers (right y-axis). Each data point (black dot) is a measured DWF value, and the vertical spread is indicative of the distribution of DWF values found across the literature, while the horizontal width of each plot is indicative of how many times, across studies, the same DWF value was found.

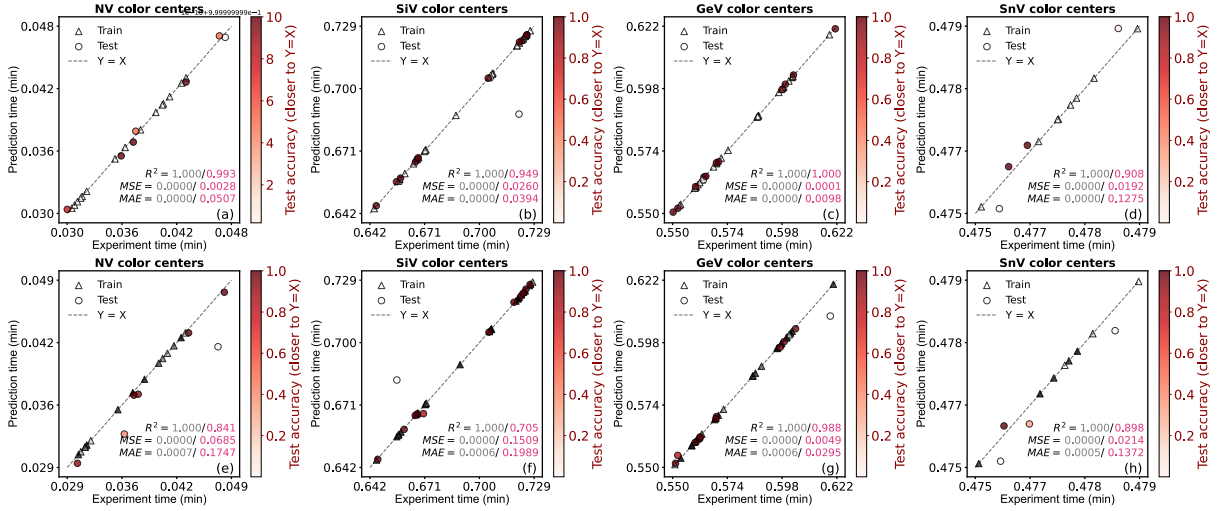


To train and test the performance of our DTR and XGB algorithms we use the same approach carried out for *Task I*. We organized the information from the literature in four databases (SI, § S2, Table ST2–ST5) and used a 70-30 split for the analysis, where 70% of the entries, at random, are used for training the DTR/XGB predictors and the remaining 30% to test them. For all the synthesis techniques, the selected target value is the DWF, whereas the input features differ for each method (SI, § S1, Table ST1). The performance of the models is based on the difference between the true experimental values and those produced by the DTR/XGB predictors (see § 3.1) and it is assessed, practically, by estimating the quantities  $R^2$ , MSE and MAE. We remark that the DWF is rarely explicitly reported in any of the experimental studies we reviewed. To determine the DWF values, we extracted the spectrum from each publication and calculated the DWF from the available data. Figure 6 is a *violin* chart that summarizes the experimental values of the DWF extracted from the literature, organized by color center and synthesis method. The vertical and horizontal spreads of each plot are indicative of the dispersion of the measured DWF values and their observation frequency, respectively.

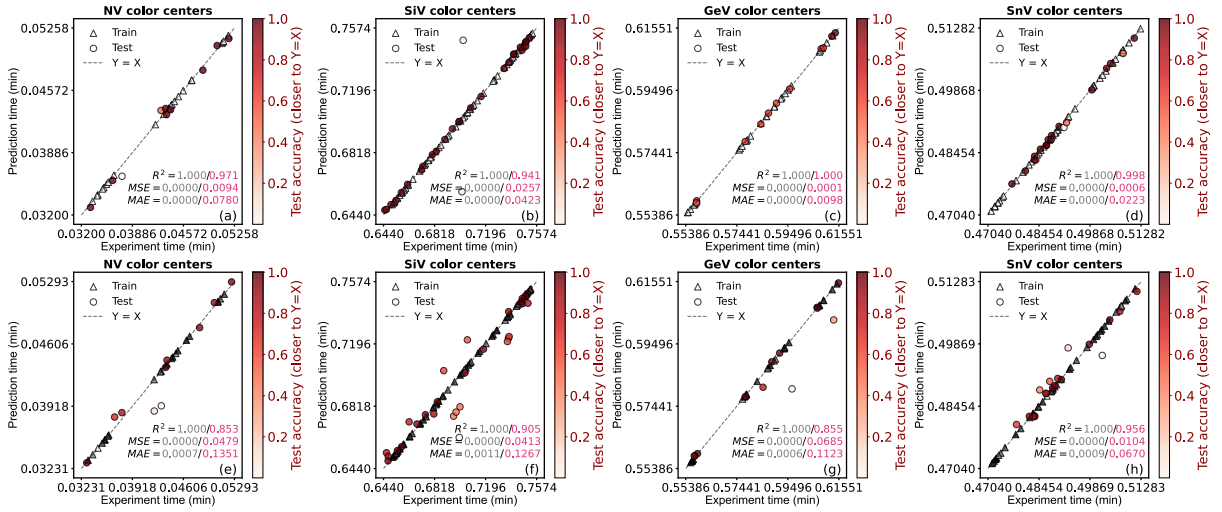
The results of the training and testing of the DTR and XGB algorithms for *Task II* are displayed in Figures 7–10 for the HPHT, CVD/MPCVD, Ion implantation and Electron/ion irradiation techniques, respectively. The estimated values for  $R^2$ , MSE and MAE are displayed in the bottom-right corner of each subpanel both for the test and training data (the latter are merely used as control); they are also summarized in the SI, § S4, Table ST8.



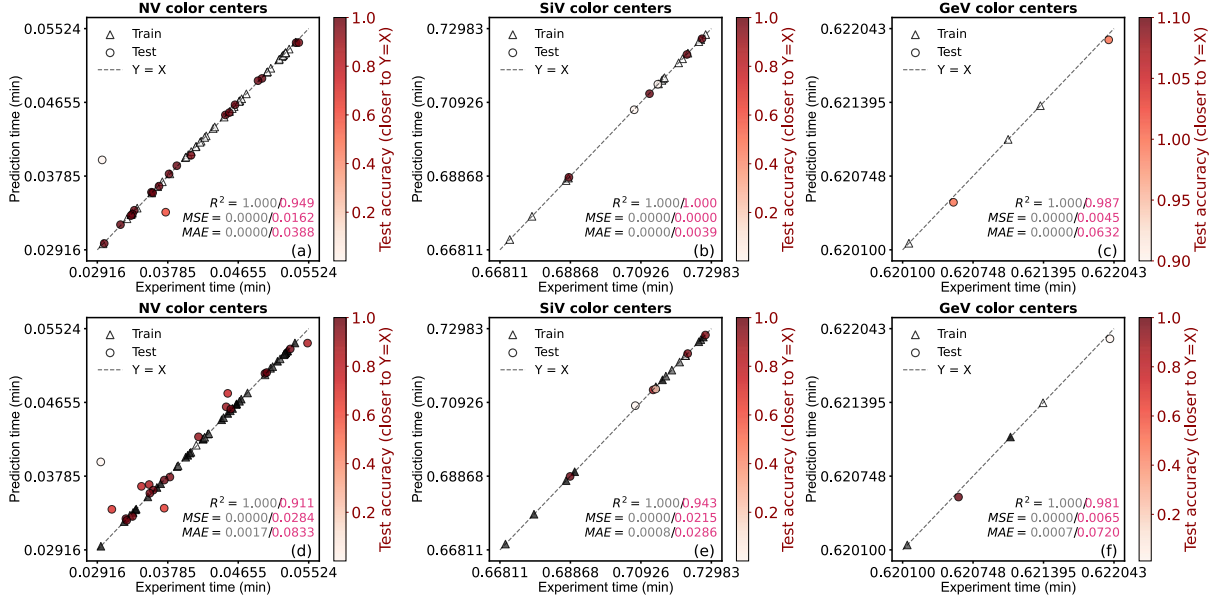
**Figure 7.** HPHT, *Task II* summary. DTR (top row) and XGB (bottom row) prediction of the DWF (target variable) from given values for pressure, temperature, run time and average size of the diamond particles (input features), organized by color centers (columns). The values for  $R^2$ , MSE and MAE for the training (triangles) and test (circles) sets are listed in the bottom-right corner of each subpanel.



**Figure 8.** CVD/MPCVD, *Task II* summary. DTR (top row) and XGB (bottom row) prediction of the DWF (target variable) from given values for chamber pressure, substrate temperature and MW plasma power (input features), organized by color centers (columns). The values for  $R^2$ , MSE and MAE for the training (triangles) and test (circles) sets are listed in the bottom-right corner of each subpanel.



**Figure 9.** Ion implantation, *Task II* summary. DTR (top row) and XGB (bottom row) prediction of the DWF (target variable) from given values for implantation energy, fluence, annealing temperature and annealing time (input features), organized by color centers (columns). The values for  $R^2$ , MSE and MAE for the training (triangles) and test (circles) sets are listed in the bottom-right corner of each subpanel.



**Figure 10.** Electron/ion irradiation, *Task II* summary. DTR (top row) and XGB (bottom row) prediction of the DWF (target variable) from given values for irradiation energy, fluence, annealing temperature and annealing time (input features), organized by color centers (columns). The values for  $R^2$ , MSE and MAE for the training (triangles) and test (circles) sets are listed in the bottom-right corner of each subpanel.

Similarly to the case of *Task I*, we can make a series of general observations related to the results of Figures 7–10 for *Task II*. Overall, the DTR and XGB models are able to predict the DWF value quite well: the estimates for  $R^2$  are consistently  $> 0.9$  and often close to  $\sim 1$ , with just a few exceptions for the XGB predictor where  $R^2 \leq 0.85$ . The value of MSE and MAE were also desirably small and consistently  $\lesssim 10\%$  of the average target value, in fact regularly within just a few percent points of it. We hypothesize that some of the larger deviations are due to our models not including categorical features such as the substrate materials or the ion source, which we excluded for practical reasons, discussed in detail in § 2.3.2. This is because a different choice of substrate or ion source (not included in our algorithms) can affect the values of other input features that appear in our models—such as annealing time or fluence—indirectly affecting the predictors’ performance. Yet, we remark that this occurs in a very small minority of cases, underscoring the overall robustness of the predictive ability of our DTR and XGB algorithms. We also found that the DTR and XGB predictors were comparable in terms of performance, with DTR slightly outperforming XGB in most cases (see explanation in § 3.1).

As a final remark, one of the merits of our meta-analysis is that it reveals that for the same color center, different techniques yield different ranges for DWF values. These are listed explicitly and in full in the SI, § S5, Table ST9, and show that ion implantation and electron/ion irradiation tend to yield DVCCs with higher DWF values—which is usually desirable for optical, photonic and quantum applications.

### 3.3. Shapley value analysis of material synthesis

In this section, we discuss the results of the Shapley value analysis for *Task I* of each synthesis technique, based on the theoretical framework introduced in § 2.3.3. Shapley value analysis is a particularly valuable tool for assessing and interpreting black-box predictive models (such as the DTR/XGB used in this study). It is model agnostic, and it helps quantify explicitly the relative interdependence between input features and target variables across all observations. This follows from the fact that the Shapley values represent the marginal contribution of individual features and of individual instances to the predicted value. It measures how much each feature,  $\phi_j^{(i)}$  in eq. (7), in every instance  $i$  (an instance here is an individual entry in the database) contributes to return a predicted value for the target variable,  $f(\mathbf{x}_i)$ , above or below the baseline value, i.e.  $\phi_0$  in eq. (7). We note that, in general, the relative weights of the features in our DTR/XGB models do not provide the complete picture. They might show that one feature is more important than the others, but not how it affects directly the prediction  $\hat{y}_i = f(\mathbf{x}_i)$  of any one particular  $i$ -th database entry. Conversely, Shapley values show exactly how much and in what direction each feature moves the prediction for an individual instance above or below the mean value,  $\phi_0$ .

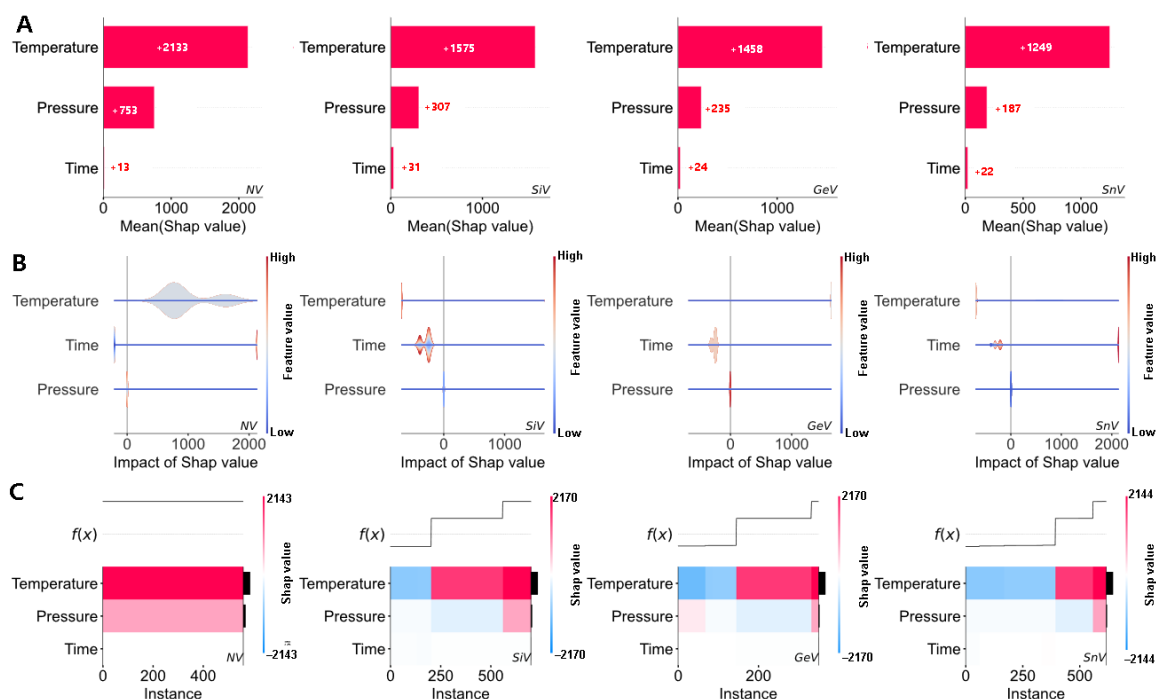
The results of the Shapley value analysis for the DTR model are shown in Figures 11–14 for the HPHT, CVD/MPCVD, Ion implantation and Electron/ion irradiation synthesis methods, respectively. Each figure is organized by color center (columns) and the results are summarized and arranged in rows of subpanels labelled A, B and C, and explained below.

Subpanels A shows the *feature importance bar* plots. In this type of plots, the various features are listed vertically, sorted from most to least important; the ranking is based on the mean absolute Shapley value for the features over all the instances. The  $x$ -axis shows the average absolute Shapley value of each feature across the entire dataset. It indicates how much, on average, each feature contributed to changing the prediction (regardless of whether it increased or decreased it). Effectively, a longer (shorter) bar indicates that a feature is more (less) influential in affecting the model’s predictions across instances. For example, Figure 11A shows that in the HPHT process for the case of the GeV center, across all the instances, the chamber temperature was significantly much more important (+1458) in causing changes in the predicted size of the synthesized diamond materials than the chamber pressure (+235); the run time rarely had any influence (+24).

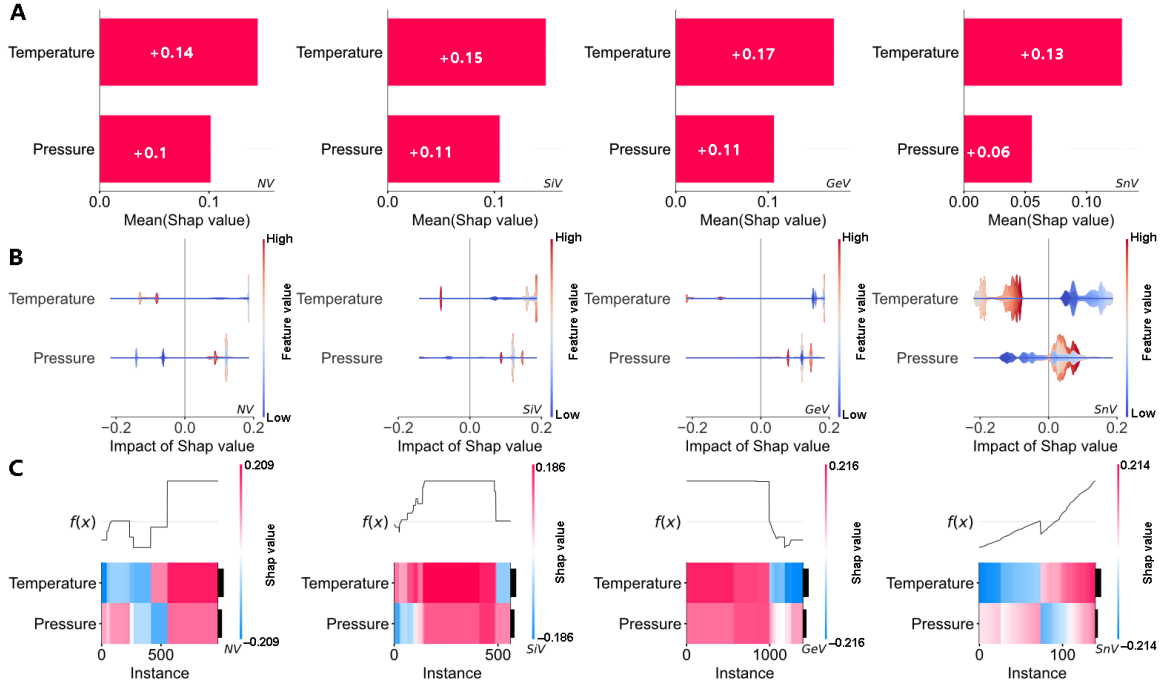
Subpanels B shows the *violin* plots. The violin plot contains several layers of information. The features are listed in order of importance from top to bottom. The lateral spread shows how variable the contribution of one feature is across observations. If a feature has positive (negative) Shapley values, it means that for those observations/instances the feature increased (decreased)

the calculated predicted value with respect to the mean value,  $\phi_0$ , of the prediction. The plot also carries information through color encoding: red (blue) colors on the positive side indicates that a high (low) feature value pushed the prediction above the average,  $\phi_0$ , while red (blue) colors on the negative side indicate that a high (low) feature value pushed the prediction below the average,  $\phi_0$ . The vertical spread is simply indicative of the density of data (more observations of the same Shapley value translate into more points stacked). We can use again the case of the GeV center in Figure 11B as an example. Here, high values (red) of the feature *Time* contributed to the output target variable *Size* by reducing (being to the left of '0') its value below the average. This occurred consistently across all instances, as shown by the fact that all the data concentrates around the same few negative values and are thus spread out vertically.

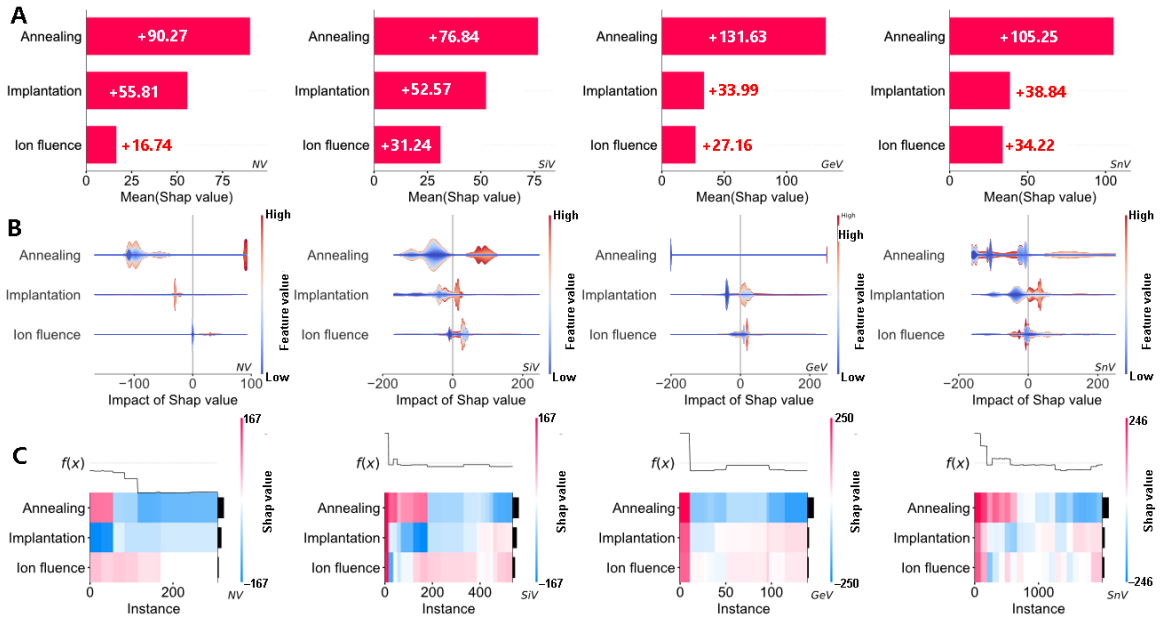
Subpanels C shows the *heatmap* plots. In these plots, the *x*-axis shows the individual instances (i.e. the entries in the database) sorted by Shapley value magnitude. The  $f(\mathbf{x})$  graph is the model prediction (eq. (6)) showing the model's output through the individual instances (again, these are sorted such that the patterns in feature contribution align with the progression of predictions). The red (blue) colors indicate a positive (negative) Shapley value that pushes the prediction above (below) the average value,  $\phi_0$ ; the white color indicates no or negligible contribution. The color changes align with the step changes in  $f(\mathbf{x})$  indicating why predictions increased or decreased. We can use again the case of the GeV center in Figure 11C as an example. The instances are ordered for convenience to show the predicted values,  $f(\mathbf{x})$ , increasing monotonically. For each value of  $f(\mathbf{x})$  the heat plot shows the positive (red), negative (blue) or neutral (white) contribution of each feature that, added, produced that value for  $f(\mathbf{x})$ .



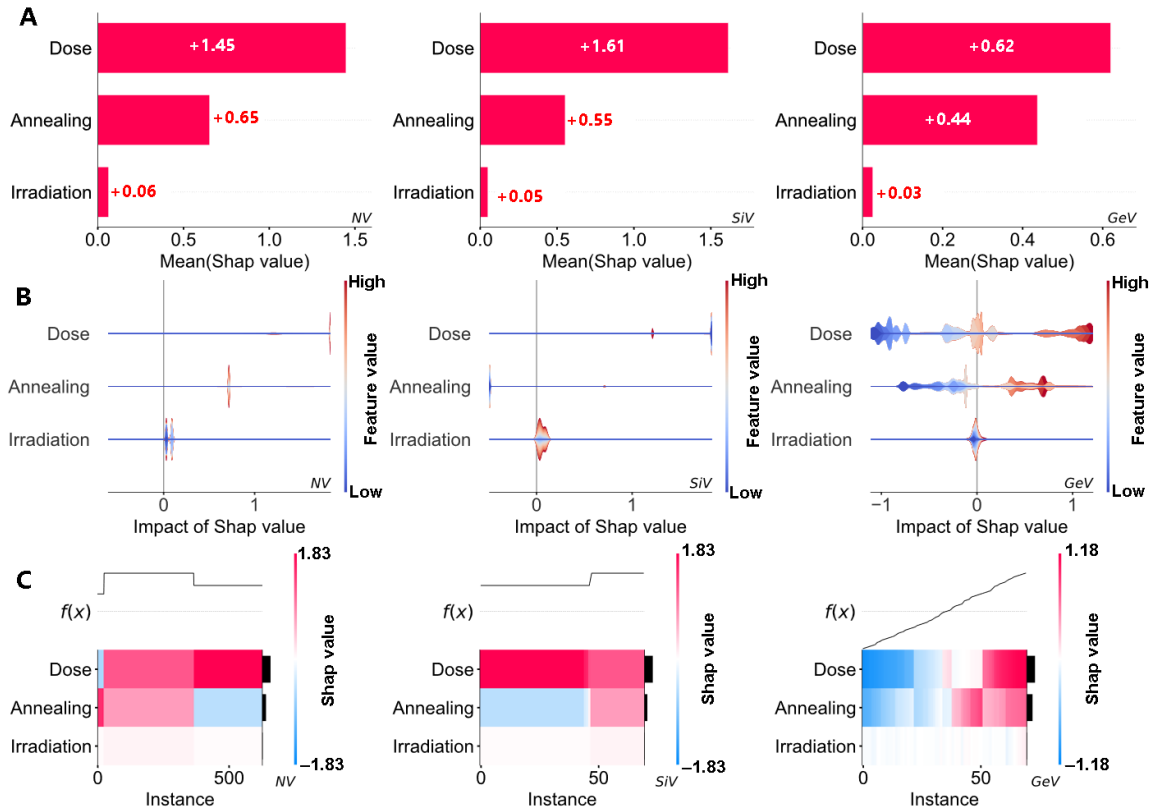
**Figure 11.** Shapley value analysis of the DTR model for HPHT synthesis. **A)** Feature importance bar plot. **B)** Violin plot. **C)** Heatmap plot. See main text for plot details.



**Figure 12.** Shapley value analysis of the DTR model for CVD/MPCVD synthesis. **A)** Feature importance bar plot. **B)** Violin plot. **C)** Heatmap plot. See main text for plot details.



**Figure 13.** Shapley value analysis of the DTR model for Ion implantation. **A)** Feature importance bar plot. **B)** Violin plot. **C)** Heatmap plot. See main text for plot details.



**Figure 14.** Shapley value analysis of the DTR model for Electron/ion irradiation. **A)** Feature importance bar plot. **B)** Violin plot. **C)** Heatmap plot. See main text for plot details.

The careful analysis of Figures 11–14 can provide detailed information about any one feature-output interdependence of interest. Here, we briefly highlight the main trends, whilst underscoring a few important subtleties associated with the analysis of Shapley values.

For HPHT synthesis, our results suggest that variations in temperature, followed by variations in pressure are the leading cause for deviations from the expected average size of the synthesized diamond materials. Conversely, time does not seem to have a significant effect.

For CVD/MPCVD, chamber pressure and substrate temperature have comparable contributions—with the latter being marginally more influential—on driving shifts in predicted MW power with respect to the expected output.

For both Ion implantation and Electron/ion irradiation, annealing temperature is the feature that mostly affects differences in annealing time. Perhaps expectedly, higher values of the ion dose and fluence correlate with higher annealing temperatures and annealing times, as—intuitively—the latter are both conducive (or necessary) to the effective repair of damage in the material caused by the irradiation process.

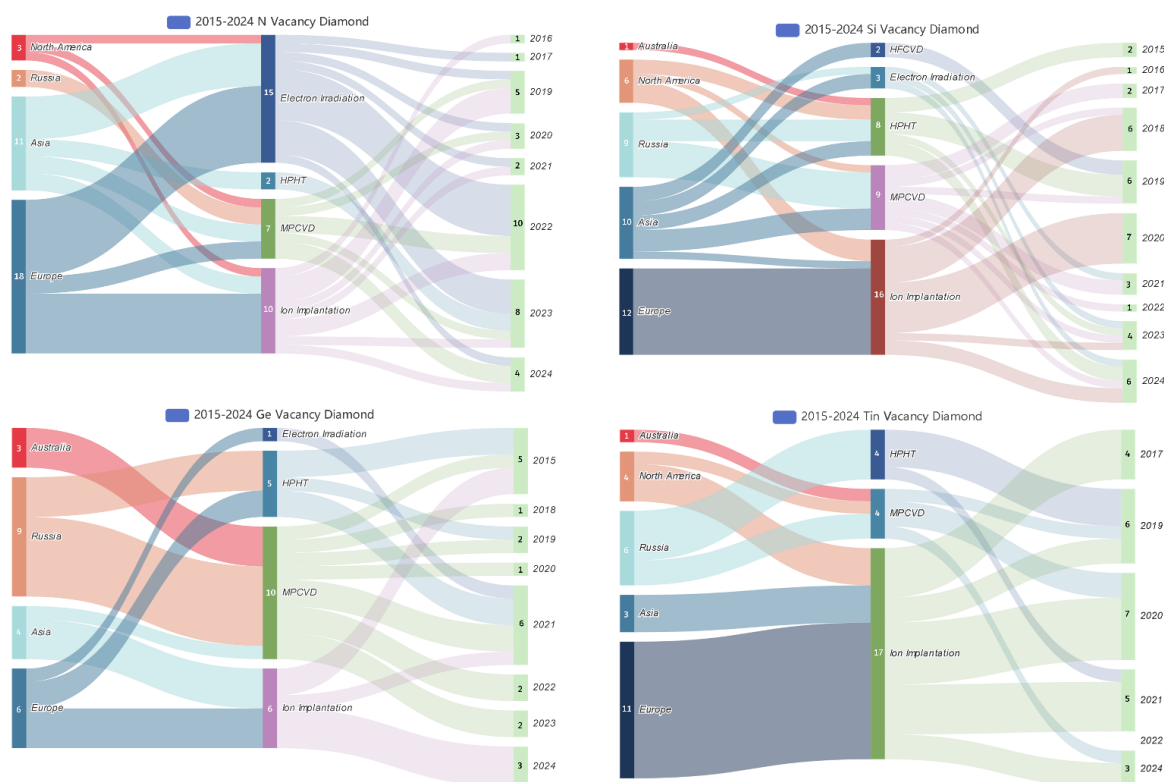
As a concluding remark to this section, we emphasize, again, that the Shapley value analysis is an instance-dependent explanation not an absolute account on relative weights in the model.



The Shapley value of the  $j$ -th feature explains for any one  $i$ -th database entry,  $\mathbf{x}_i$ , how much this  $j$ -th feature shifts the prediction,  $\hat{y}_i$ , away from the baseline (expected mean output,  $\phi_0$ ). These shifts calculated for all instances in the database can then be aggregated, for example, in the Shapley feature importance bar plot, to make general inferences. Therefore, theoretically, if a feature did not vary much in the dataset, its Shapley value would be small, regardless of how large or small its weight is in the model, and vice versa. The Shapley analysis measures the impact on the prediction, not on the model parameters.

#### 4. Worldwide trends in diamond synthesis

The database we compiled is rich in information that goes beyond process parameters and fabrication protocols. In this section, we leverage this information to produce Sankey diagrams that visualize some of the main trends and practices in diamond synthesis over the last decade—as inferred from the studies we reviewed in our meta-analysis. The diagram (Fig. 15) has three main nodes: geographical region, synthesis method and year. Specific trends of interest can be deduced directly from the analysis of the diagram; below we highlight a few main notable points.



**Figure 15.** Sankey diagrams for the synthesis of diamond materials containing NV, SiV, GeV, and SnV color centers, organized by geographical region, synthesis method, and year. The diagrams were extracted from a subset of studies surveyed for our meta-analysis, focusing on data from the last decade.

Geographically, most of the studies we gathered were conducted by groups in Europe (47/119,



39%), Asia (28/119, 23%) and Russia (26/119, 22%). Regarding diamond synthesis and DVCC inclusion, ion implantation and electron/ion irradiation are, across the board, widely used for color centers incorporation in diamond, with many of the studies surveyed coming out of Europe (37/68, 54%) and Asia (18/68, 26%). The number of studies we collated was split roughly evenly with respect to the synthesis method, with a slight edge for CVD/MPCVD (30/49, 61%) over HPHT (19/49, 39%). Electron irradiation is more commonly used for incorporation of NV color centers (15/25, 60%), whereas other color centers are mostly created through ion implantation (23/27, 85%). This is expected because nitrogen is abundant in diamonds, and electron irradiation is usually sufficient to create vacancies and—usually upon annealing—stable NV centers. Conversely, for SiV, GeV, and SnV, the incorporation of the relevant dopant via implantation is either needed or desirable, especially in studies where it is performed deterministically (e.g., to leverage emission from single color centers). We note that MPCVD is also commonly used for incorporation of color centers in diamond directly during growth, especially for GeV (10/22, 45%) and SiV (9/38, 24%). For SnV the preferred incorporation method is ion implantation (17/25, 68%).

## 5. Conclusions

In this work we performed a thorough review and meta-analysis of studies involving the synthesis of diamond materials via HPHT, CVD/MPCVD, Ion implantation, and Electron/ion irradiation, over the last decade. Data from approximately 60 published articles was organized into four databases, including a total of 170 data sets and 1692 entries. These were used to train two ML algorithms: DTR and XGB, with the goal of using these models to select diamond synthesis techniques and parameters that would yield materials and DVCCs with the desired characteristics, properties and quality. We organized our analysis in two streams: one aimed at optimizing diamond material synthesis (*Task I*), and one aimed at optimizing the yield of high-quality color centers in the fabricated diamond (*Task II*). We characterized the predictive performance of our DTR/XGB models using traditional statistical indicators, namely  $R^2$ , MSE, and MAE, and performed a Shapley value analysis to complement their interpretation.

Overall, our results show that the proposed DTR and XGB models are robust predictive algorithms that can leverage knowledge of input features (i.e. synthesis parameters) to make accurate predictions about the specific properties and quality (target variables) of synthesized diamond materials. These models can also be effectively used in reverse to tweak and control—optimally and efficiently—the value of specific synthesis parameters to yield materials with ad hoc characteristics. Thus, the data-driven framework we built is a powerful addition to the tools available to researchers and material scientists involved in the synthesis of diamond materials and diamond color centers.

## 6. Supporting Information

Input features and target variable used in training/testing the algorithms (§ S1, Table ST1), and schematic representation of XGB (§ S1, Fig. SF1). Databases used to train the DTR/XGB algorithms (§ S2, Tables ST2–ST5). Hyperparameters and their value ranges used in the models (§ S3, Table ST6). Summary of the statistical estimators we used to measure the performance of our models for *Tasks I* and *II* (§ S4, Tables ST7 and ST8). Summary of the ranges of DWF values organized by color center and synthesis technique (§ S5, Table ST9).

## 7. Contributions and Acknowledgements

Collection and meta-analysis of the data was performed by Z. J. Theoretical modeling and data analysis were performed by Z.J., C.B and G.G. The analysis and interpretability of the data rationale were provided by M.P and C.B. All authors contributed to writing the manuscript. The authors declare that they have no conflicts of interest or competing interests related to the research presented in this study.

The Horizon Europe (HORIZON), HORIZON-MSCA-2023-PF-01, Grant101149632–GDSNL, project UID 00481 Centre for Mechanical Technology and Automation (TEMA) and project CarboNCT, 2022.03596.PTDC (DOI: 10.54499/2022.03596.PTDC) are acknowledged.

The Natural Sciences and Engineering Research Council of Canada (DGEER-2021-00234) and the Canada Foundation for Innovation (John R. Evans Leaders Fund #41173) are acknowledged.

## 8. Conflict of Interest

The authors declare no conflict of interest.

## 9. Data Availability

The Supporting Information file contains additional details about the databases, models and their performance. The data and code of the DTR and XGB algorithms used in this study are available in a public repository (linked below).

<https://github.com/howgoods/Diamond-color-centers-python-code>

## References

- [1] T. Baba, M. Iizawa, K. Takenaka, K. Kimura, A. Kawasaki, T. Taniguchi, M. Miyakawa, H. Okazaki, O. Hanaizumi, S. Onoda, *physica status solidi (a)* **2024**, *n/a*, 2400303.
- [2] H. Bernien, B. Hensen, W. Pfaff, G. Koolstra, M. S. Blok, L. Robledo, T. H. Taminiau, M. Markham, D. J. Twitchen, L. Childress, R. Hanson, *Nature* **2013**, *497*, 86.
- [3] W. J. Munro, K. Azuma, K. Tamaki, K. Nemoto, *IEEE Journal of Selected Topics in Quantum Electronics* **2015**, *21*, 78.
- [4] R. H. H. Bernien A. E. Dréau, A. Reiserer, N. Kalb, M. S. Blok, J. Ruitenber, R. F. L. Vermeulen, R. N. Schouten, C. Abellán, W. Amaya, V. Pruneri, M. W. Mitchell, M. Markham, D. J. Twitchen, D. Elkouss, S. Wehner, T. H. Taminiau, *Quantum Physics* **2015**, *526*, 682.
- [5] W. Pfaff, B. J. Hensen, H. Bernien, S. B. van Dam, M. S. Blok, T. H. Taminiau, M. J. Tiggelman, R. N. Schouten, M. Markham, D. J. Twitchen, R. Hanson, *Science* **2014**, *345*, 532.
- [6] J. R. Weber, W. F. Koehl, J. B. Varley, A. Janotti, B. B. Buckley, C. G. V. de Walle, D. D. Awschalom, *Proceedings of the National Academy of Sciences* **2010**, *107*, 8513.
- [7] B. Gardas, J. Dziarmaga, W. H. Zurek, M. Zwolak, *Scientific Reports* **2018**, *8*, 4539.
- [8] R. H. Lavroff, D. L. Pennington, A. S. Hua, B. Y. Li, J. A. Williams, A. N. Alexandrova, *J. Phys. Chem. A* **2021**, *125*, 9567.
- [9] I. Buluta, F. Nori, *Science* **2009**, *326*, 108.
- [10] J. Cai, A. Retzker, F. Jelezko, M. B. Plenio, *Nature Physics* **2013**, *9*, 168.
- [11] D. D. Awschalom, R. Hanson, J. Wrachtrup, B. B. Zhou, *Nature Photonics* **2018**, *12*, 516.
- [12] C. L. Degen, F. Reinhard, P. Cappellaro, *Rev. Mod. Phys.* **2017**, *89*, 035002.
- [13] S. Vaidya, X. Gao, S. Dikshit, I. Aharonovich, T. Li, *Advances in Physics: X* **2023**, *8*, DOI 10.1080/23746149.2023.2206049.
- [14] B. Sanchez-Lengeling, A. Aspuru-Guzik, *Science* **2018**, *361*, 360.
- [15] C. Sutton, M. Boley, L. M. Ghiringhelli, M. Rupp, J. Vreeken, M. Scheffler, *Nature Communications* **2020**, *11*, 4428.
- [16] H. Moon, D. T. Lennon, J. Kirkpatrick, N. M. van Esbroeck, L. C. Camenzind, L. Yu, F. Vigneau, D. M. Zumbühl, G. A. D. Briggs, M. A. Osborne, D. Sejdinovic, E. A. Laird, N. Ares, *Nature Communications* **2020**, *11*, 4161.
- [17] S. Bhowmik, A. G. Rajan, *iScience* **2022**, *25*, 103832.
- [18] B. Ryu, L. Wang, H. Pu, M. K. Y. Chan, J. Chen, *Chem. Soc. Rev.* **2022**, *51*, 1899.
- [19] H. Tao, T. Wu, M. Aldeghi, T. C. Wu, A. Aspuru-Guzik, E. Kumacheva, *Nature Reviews Materials* **2021**, *6*, 701.
- [20] S. A. Tawfik, O. Isayev, C. Stampfl, J. Shapter, D. A. Winkler, M. J. Ford, *Advanced Theory and Simulations* **2019**, *2*, 1800128.
- [21] C. H. Chan, M. Sun, B. Huang, *EcoMat* **2022**, *4*, e12194.
- [22] N. M. van Esbroeck, D. T. Lennon, H. Moon, V. Nguyen, F. Vigneau, L. C. Camenzind, L. Yu, D. M. Zumbühl, G. A. D. Briggs, D. Sejdinovic, N. Ares, *New Journal of Physics* **2020**, *22*, 095003.
- [23] K. Kojima, M. H. Tahersima, T. Koike-Akino, D. Jha, Y. Tang, Y. Wang, K. Parsons, F. Sang, J. Klamkin, in *2020 Optical Fiber Communications Conference and Exhibition*

- (OFC), **2020**, pp. 1–3.
- [24] J. Jiang, M. Chen, J. A. Fan, *Nature Reviews Materials* **2021**, 6, 679.
  - [25] E. Chávez-Angel, M. B. Eriksen, A. Castro-Alvarez, J. H. Garcia, M. Botifoll, O. Avalos-Ovando, J. Arbiol, A. Mugarza, *Advanced Intelligent Systems* **2025**, 2400986.
  - [26] M. W. Doherty, N. B. Manson, P. Delaney, F. Jelezko, J. Wrachtrup, L. C. L. Hollenberg, *Physics Reports* **2013**, 528, 1.
  - [27] C. Bradac, W. Gao, J. Forneris, M. E. Trusheim, I. Aharonovich, *Nature Communications* **2019**, 10, 5625.
  - [28] H. Morishita, T. Tashima, D. Mima, H. Kato, T. Makino, S. Yamasaki, M. Fujiwara, N. Mizuochi, *Scientific Reports* **2019**, 9, 13318.
  - [29] M. W. N. Ngambou, C. Pellet-Mary, O. Brinza, A. Antonino, G. Hetet, A. Tallaire, F. Bénédict, J. Achard, *Diamond and Related Materials* **2022**, 123, 108884.
  - [30] X. Gao, C. Yu, S. Zhang, H. Lin, J. Guo, M. Ma, Z. Feng, F.-W. Sun, *Diamond and Related Materials* **2023**, 139, 110348.
  - [31] T. Liu, F. Lyu, T. Shao, D. Zou, W. Shen, Y. Guo, Y. Zhong, C. Chen, L. Yi, Z. Zhang, A. H. Shen, *Materials Futures* **2024**, 3, 045701.
  - [32] Z. Jiang, E. Wang, J. Ying, C. Zhang, J. Du, G. Wang, Y. Pang, *IEEE Transactions on Electron Devices* **2024**, 71, 6236.
  - [33] A. F. L. Poulsen, J. D. Clement, J. L. Webb, R. H. Jensen, L. Troise, K. Berg-Sørensen, A. Huck, U. L. Andersen, *Physical Review B* **2022**, 106, 014202.
  - [34] K. Choudhary, B. DeCost, C. Chen, A. Jain, F. Tavazza, R. Cohn, C. W. Park, A. Choudhary, A. Agrawal, S. J. L. Billinge, E. Holm, S. P. Ong, C. Wolverton, *npj Computational Materials* **2022**, 8, 59.
  - [35] S. W. Im, D. Zhang, J. H. Han, R. M. Kim, C. Choi, Y. M. Kim, K. T. Nam, *Nature Materials* **2024**, 23, 977.
  - [36] F. P. BUNDY, H. T. HALL, H. M. STRONG, R. H. WENTORFJUN., *Nature* **1955**, 176, 51.
  - [37] B. B. Pate, *Surface Science* **1986**, 165, 83.
  - [38] X. Zhang, K.-Y. Liu, F. Li, X. Liu, S. Duan, J.-N. Wang, G.-Q. Liu, X.-Y. Pan, X. Chen, P. Zhang, Y. Ma, C. Chen, *Advanced Functional Materials* **2023**, 33, 2309586.
  - [39] P. Bhattacharyya, W. Chen, X. Huang, S. Chatterjee, B. Huang, B. Kobrin, Y. Lyu, T. J. Smart, M. Block, E. Wang, Z. Wang, W. Wu, S. Hsieh, H. Ma, S. Mandyam, B. Chen, E. Davis, Z. M. Geballe, C. Zu, V. Struzhkin, R. Jeanloz, J. E. Moore, T. Cui, G. Galli, B. I. Halperin, C. R. Laumann, N. Y. Yao, *Nature* **2024**, 627, 73.
  - [40] R. N. Patel, T. Schröder, N. Wan, L. Li, S. L. Mouradian, E. H. Chen, D. R. Englund, *Light: Science & Applications* **2016**, 5, e16032.
  - [41] T. Luo, L. Lindner, J. Langer, V. Cimalla, X. Vidal, F. Hahl, C. Schreyvogel, S. Onoda, S. Ishii, T. Ohshima, D. Wang, D. A. Simpson, B. C. Johnson, M. Capelli, R. Blinder, J. Jeske, *New Journal of Physics* **2022**, 24, 033030.
  - [42] F. A. Hahl, L. Lindner, X. Vidal, T. Luo, T. Ohshima, S. Onoda, S. Ishii, A. M. Zaitsev, M. Capelli, B. C. Gibson, A. D. Greentree, J. Jeske, *Science Advances* **n.d.**, 8, eabn7192.
  - [43] L. Lindner, F. A. Hahl, T. Luo, G. N. Antonio, X. Vidal, M. Rattunde, T. Ohshima, J. Sacher, Q. Sun, M. Capelli, B. C. Gibson, A. D. Greentree, R. Quay, J. Jeske, *Science*

- Advances* **2024**, *10*, eadj3933.
- [44] S. Zhang, Z. Ye, Y. Zhu, M. Ye, S. Ye, M. He, S. Li, S. Cai, C. Zhao, Y. Hang, *Crystal Growth & Design* **2024**, *24*, 6701.
  - [45] W. G. Eversole, *Synthesis of Diamond*, **1962**, US3030188A, US3030187A.
  - [46] V. V. S. S. Srikanth, *Proceedings of the Institution of Mechanical Engineers, Part C* **2012**, *226*, 303.
  - [47] Y. S. Nubuo Setaka, *Journal of Materials Science* **1981**, *16*, 1728.
  - [48] S. Matsumoto, Y. Sato, M. Kamo, N. Setaka, *Japanese Journal of Applied Physics* **1982**, *21*, L183.
  - [49] S. Matsumoto, Y. Sato, M. Tsutsumi, N. Setaka, *Journal of Materials Science* **1982**, *17*, 3106.
  - [50] M. Kamo, Y. Sato, S. Matsumoto, N. Setaka, *Journal of Crystal Growth* **1983**, *62*, 642.
  - [51] S. Matsumoto, Y. Matsui, *Journal of Materials Science* **1983**, *18*, 1785.
  - [52] P. K. Bachmann, *Materials Science and Engineering: B* **1991**, *7*, 231.
  - [53] Y. M. H. Yusuke Tominaga, *Diamond & Related Materials* **2023**, *140*, DOI 10.1016/j.diamond.2023.110543.
  - [54] A. R. H. H. C. Thejaswini, *Diamond & Related Materials* **2014**, *48*, 24.
  - [55] S. Malykhin, Y. Mindarava, R. Ismagilov, A. Orekhov, F. Jelezko, A. Obraztsov, *physica status solidi (b)* **2019**, *256*, 1800721.
  - [56] Y.-K. Tzeng, J. L. Zhang, H. Lu, H. Ishiwata, J. Dahl, R. M. K. Carlson, H. Yan, P. R. Schreiner, J. Vučković, Z.-X. Shen, N. Melosh, S. Chu, *Nano Lett.* **2017**, *17*, 1489.
  - [57] K. N. Boldyrev, V. S. Sedov, D. E. P. Vanpoucke, V. G. Ralchenko, B. N. Mavrin, *Diamond and Related Materials* **2022**, *126*, 109049.
  - [58] V. Sedov, A. Martyanov, A. Neliubov, I. Tiazhelov, S. Savin, I. Eremchev, M. Eremchev, M. Pavlenko, S. Mandal, V. Ralchenko, A. Naumov, *Philosophical Transactions A* **2024**, *382*, 20230167.
  - [59] J. C. Mendes, M. Liehr, C. Li, *Materials* **2022**, *15*, DOI 10.3390/ma15020415.
  - [60] S. Praver, *Diamond and Related Materials* **1995**, *4*, 862.
  - [61] T. Herzig, T. Lüthmann, P. Räcke, C. Scheuner, S. Pezzagna, J. Meijer, in *Diamond for Quantum Applications Part 2* (Eds.: C. E. Nebel, I. Aharonovich, N. Mizuochi, M. Hatano), Elsevier, **2021**, pp. 1–30.
  - [62] M. Vićentijević, M. Jakšić, T. Suligoj, *Scientific Reports* **2023**, *13*, 13483.
  - [63] K. Groot-Berning, G. Jacob, C. Osterkamp, F. Jelezko, F. Schmidt-Kaler, *New Journal of Physics* **2021**, *23*, 063067.
  - [64] T. Iwasaki, F. Ishibashi, M. Miyamoto, Y. Doi, S. Kobayashi, T. Miyazaki, K. Tahara, K. D. Jahnke, L. J. Rogers, B. Naydenov, F. Jelezko, S. Yamasaki, S. Nagamachi, T. Inubushi, N. Mizuochi, M. Hatano, *Sci. Rep.* **2015**, *5*, 12882.
  - [65] F. Fávaro de Oliveira, S. A. Momenzadeh, D. Antonov, J. Scharpf, C. Osterkamp, B. Naydenov, F. Jelezko, A. Denisenko, J. Wrachtrup, *Nano Lett.* **2016**, *16*, 2228.
  - [66] T. Schröder, M. E. Trusheim, M. Walsh, L. Li, J. Zheng, M. Schukraft, A. Sipahigil, R. E. Evans, D. D. Sukachev, C. T. Nguyen, J. L. Pacheco, R. M. Camacho, E. S. Bielejec, M. D. Lukin, D. Englund, *Nature Communications* **2017**, *8*, 15376.
  - [67] S. D. Tchernij, T. Herzig, J. Forneris, J. Küpper, S. Pezzagna, P. Traina, E. Moreva, I. P.

- Degiovanni, G. Brida, N. Skukan, M. Genovese, M. Jakšić, J. Meijer, P. Olivero, *ACS Photonics* **2017**, *4*, 2580.
- [68] R. A. Parker, J. Arjona Martínez, K. C. Chen, A. M. Stramma, I. B. Harris, C. P. Michaels, M. E. Trusheim, M. Hayhurst Appel, C. M. Purser, W. G. Roth, D. Englund, M. Atatüre, *Nature Photonics* **2024**, *18*, 156.
- [69] U. Wahl, J. G. Correia, Â. Costa, A. Lamelas, V. Amaral, K. Johnston, G. Magchiels, S. M. Tunhuma, A. Vantomme, L. M. C. Pereira, *Materials for Quantum Technology* **2024**, *4*, 025101.
- [70] V. Yurgens, A. Corazza, J. A. Zuber, M. Gruet, M. Kasperczyk, B. J. Shields, R. J. Warburton, Y. Fontana, P. Maletinsky, *Applied Physics Letters* **2022**, *121*, 234001.
- [71] J. Lang, S. Häußler, J. Fuhrmann, R. Waltrich, S. Laddha, J. Scharpf, A. Kubanek, B. Naydenov, F. Jelezko, *Applied Physics Letters* **2020**, *116*, 064001.
- [72] E. A. Ekimov, S. G. Lyapin, K. N. Boldyrev, M. V. Kondrin, R. Khmelnitskiy, V. A. Gavva, T. V. Kotereva, M. N. Popova, *JETP Letters* **2015**, *102*, 701.
- [73] R. Debroux, C. P. Michaels, C. M. Purser, N. Wan, M. E. Trusheim, J. Arjona Martínez, R. A. Parker, A. M. Stramma, K. C. Chen, L. de Santis, E. M. Alexeev, A. C. Ferrari, D. Englund, D. A. Gangloff, M. Atatüre, *Phys. Rev. X* **2021**, *11*, 041041.
- [74] T. Teraji, C. Shinei, Y. Masuyama, M. Miyakawa, T. Taniguchi, *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences* **2023**, *382*, 20220322.
- [75] E. A. Ekimov, M. V. Kondrin, V. S. Krivobok, A. A. Khomich, I. I. Vlasov, R. A. Khmelnitskiy, T. Iwasaki, M. Hatano, *Diamond and Related Materials* **2019**, *93*, 75.
- [76] J. R. Rabeau, P. Reichart, G. Tamanyan, D. N. Jamieson, S. Prawer, F. Jelezko, T. Gaebel, I. Popa, M. Domhan, J. Wrachtrup, *Applied Physics Letters* **2006**, *88*, DOI 10.1063/1.2158700.
- [77] J. O. Orwa, C. Santori, K. M. C. Fu, B. Gibson, D. Simpson, I. Aharonovich, A. Stacey, A. Cimmino, P. Balog, M. Markham, D. Twitchen, A. D. Greentree, R. G. Beausoleil, S. Prawer, *Journal of Applied Physics* **2011**, *109*, 083530.
- [78] S. Pezzagna, B. Naydenov, F. Jelezko, J. Wrachtrup, J. Meijer, *New Journal of Physics* **2010**, *12*, 065017.
- [79] Y. Mita, *Phys. Rev. B* **1996**, *53*, 11360.
- [80] J. Meijer, T. Vogel, B. Burchard, I. W. Rangelow, L. Bischoff, J. Wrachtrup, M. Domhan, F. Jelezko, W. Schnitzler, S. A. Schulz, K. Singer, F. Schmidt-Kaler, *Applied Physics A* **2006**, *83*, 321.
- [81] E. Fraczek, V. G. Savitski, M. Dale, B. G. Breeze, P. Diggle, M. Markham, A. Bennett, H. Dhillon, M. E. Newton, A. J. Kemp, *Opt. Mater. Express* **2017**, *7*, 2571.
- [82] C. Yang, Z. Mi, H. Jin, T. Venkatesan, R. Vispute, A. A. Bettiol, *Carbon* **2023**, *203*, 842.
- [83] A. E. Rugar, C. Dory, S. Sun, J. Vučković, *Phys. Rev. B* **2019**, *99*, 205417.
- [84] T. G. I. van Wijk, E. A. Melan, R. M. Joy, E. Y. Guillaume, P. Pobedinskas, K. Haenen, D. E. P. Vanpoucke, *Carbon* **2025**, *234*, 119928.
- [85] A. A. Razgulov, S. G. Lyapin, A. P. Novikov, E. A. Ekimov, *Phys. Rev. B* **2025**, *111*, 024114.
- [86] J. Görlitz, D. Herrmann, P. Fuchs, T. Iwasaki, T. Taniguchi, D. Rogalla, D. Hardeman, P.-

- O. Colard, M. Markham, M. Hatano, C. Becher, *npj Quantum Information* **2022**, 8, 45.
- [87] I. B. W. Harris, D. Englund, *Phys. Rev. B* **2024**, 109, 085414.
- [88] Y.-K. Tzeng, F. Ke, C. Jia, Y. Liu, S. Park, M. Han, M. Frost, X. Cai, W. L. Mao, R. C. Ewing, Y. Cui, T. P. Devereaux, Y. Lin, S. Chu, *Nature Communications* **2024**, 15, 7251.
- [89] J. M. Smith, S. A. Meynell, A. C. B. Jayich, J. Meijer, *Nanophotonics* **2019**, 8, 1889.
- [90] M. Krzywinski, N. Altman, *Nature Methods* **2017**, 14, 757.
- [91] Jerome H. Friedman, *The Annals of Statistics* **2001**, 29, 1189.
- [92] L. S. Shapley, in (Eds.: H. W. Kuhn, A. W. Tucker), Princeton University Press, **1953**, pp. 307–318.
- [93] J. Nemecek, T. Pevny, J. Marecek, **2024**.
- [94] A. Althnian, D. AlSaeed, H. Al-Baity, A. Samha, A. B. Dris, N. Alzakari, A. Abou Elwafa, H. Kurdi, *Applied Sciences* **2021**, 11, DOI 10.3390/app11020796.
- [95] J Walker, *Reports on Progress in Physics* **1979**, 42, 1605.
- [96] E. Neu, D. Steinmetz, J. Riedrich-Möller, S. Gsell, M. Fischer, M. Schreck, C. Becher, *New Journal of Physics* **2011**, 13, 025012.
- [97] G. Thiering, A. Gali, in *Semiconductors and Semimetals* (Eds.: C. E. Nebel, I. Aharonovich, N. Mizuochi, M. Hatano), Elsevier, **2020**, pp. 1–36.
- [98] M. Radulaski, J. L. Zhang, Y.-K. Tzeng, K. G. Lagoudakis, H. Ishiwata, C. Dory, K. A. Fischer, Y. A. Kelaita, S. Sun, P. C. Maurer, K. Alassaad, G. Ferro, Z.-X. Shen, N. A. Melosh, S. Chu, J. Vučković, *Laser & Photonics Reviews* **2019**, 13, 1800316.

## Supplementary Information

### Prediction of synthesis parameters for N, Si, Ge and Sn diamond vacancy centers using machine learning

Zhi Jiang,<sup>1</sup> Marco Peres,<sup>2,3,4</sup> Carlo Bradac,<sup>5</sup> and Gil Gonçalves<sup>1\*</sup>

1. TEMA, Department of Mechanical Engineering, University of Aveiro, Aveiro 3810-193 Portugal
2. Instituto de Engenharia de Sistemas e Computadores - Microssistemas e Nanotecnologia, Rua Alves Redol 9, 1000-029 Lisboa, Portugal
3. IPFN, Instituto Superior Técnico, University of Lisbon, Av. Rovisco Pais 1, 1049 001 Lisbon, Portugal
4. DECN, Instituto Superior Técnico, University of Lisbon, Estrada Nacional 10 (km 139.7), 2695 066 Bobadela, Portugal
5. Department of Physics and Astronomy, Trent University, 1600 West Bank Drive, Peterborough, ON, K9L 0G2, Canada

\* Corresponding author: Gil Gonçalves, [ggoncalves@ua.pt](mailto:ggoncalves@ua.pt)

#### S1. Task features and target values; XGB algorithm

Table ST1 below summarizes the input features and target variables for each diamond synthesis method.

**Table ST1.** Synthesis input features and target variables

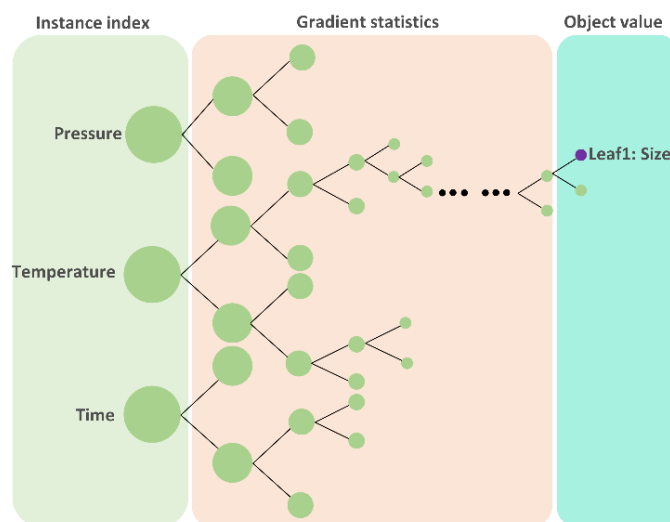
	<i>Task I</i>		<i>Task II</i>	
	Target Variable	Input Features	Target Variable	Input Features
<b>HPHT</b>	Avg. particle diameter	Chamber pressure Chamber temperature Run time	DWF of the DVCCs	Chamber pressure Chamber temperature Run time Avg. particle diameter
<b>CVD/MPCVD</b>	MW power	Chamber pressure Substrate temperature	DWF of the DVCCs	Chamber pressure Substrate temperature MW plasma power
<b>Ion implantation</b>	Annealing time	Implantation energy Fluence Annealing temperature	DWF of the DVCCs	Implantation energy Fluence Annealing temperature Annealing time
<b>Electron/ion irradiation</b>	Annealing time	Irradiation energy Fluence Annealing temperature	DWF of the DVCCs	Irradiation energy Fluence Annealing temperature Annealing time

**Table ST1.** List of target variables and input features organized by *Task* and synthesis method (HPHT, CVD/MPCVD, Ion implantation, Electron/ion irradiation).

Figure SF1 illustrates schematically the application of the XGB algorithm to the synthesis of diamond material via HPHT, used here as an example. Pressure, temperature and run time of every different row in the HPHT database (SI, § S2, Table ST2) are set as the inputs Fig. SF1, input layer – green left panel). The XGB (or



DTR) algorithm produces and optimizes the predictor (Fig. SF1, middle layer – orange middle panel). For a new, unknown set of pressure, temperature and run time input values, the optimized predictor then makes a prediction (main text, eq. (2) for DTR and eq. (5) for XGB) for the average size of the synthesized diamond material, which is the output of interest (Fig. SF1, output layer – cyan right panel).



**Figure SF1.** Schematic representation of the Extreme Gradient Boosting (XGB) algorithm for HPHT diamond synthesis.

## S2. Database

To train/test our algorithms we compiled a database from 60 experimental studies. The database has 170 data-sets and 1692 entries organized in Tables ST1–ST4. The tables are arranged by type of color center and show the values and specifications for the critical synthesis parameters discussed in the main text (§§ 2.1.1–2.1.4).

**Table ST2.** HPHT sub-database

Year	Carbon source	Vacancy source	Color centers	Pressure range (GPa)	Temperature range (°C)	Time (min)	Size (μm)	References
2023	Graphite	TiN, Fe-Co-Ti	NV	5.6 – 6.0	1645 – 1655	3595 – 3605	398 – 402	[78]
2023	Graphite	NaN <sub>2</sub> /Ba(N <sub>3</sub> ) <sub>2</sub> , Fe-Al-Ti	NV	5.2 – 5.8	1345 – 1355	4795 – 4805	492 – 502	[38–44]
2019	C <sub>10</sub> H <sub>16</sub>	C <sub>12</sub> H <sub>36</sub> Si <sub>5</sub> +C <sub>12</sub> H <sub>20</sub> Si	SiV	9.3 – 9.7	895 – 905	0.43 – 0.47	5995 – 6005	[79]
				8.3 – 8.7	1795 – 1805	58 – 62	29 – 31	
2019	CH <sub>2</sub> O+HCl+ C <sub>2</sub> H <sub>5</sub> N	Si(OC <sub>2</sub> H <sub>5</sub> ) <sub>4</sub>	SiV	18 – 22	1495 – 1505		0.08 – 0.12	[96]
				23 – 27	1695 – 1705	28 – 32	0.18 – 0.22	
2019	C <sub>10</sub> H <sub>8</sub>	C <sub>12</sub> H <sub>36</sub> Si <sub>5</sub>	SiV	7.8 – 8.2	1295 – 1305	0.1 – 0.2	0.01 – 0.12	[97]
2023	C <sub>10</sub> H <sub>8</sub> +CF <sub>4</sub> 1.1	C <sub>12</sub> H <sub>36</sub> Si <sub>5</sub>	SiV	7.8 – 8.2	1395 – 1405	0.18 – 0.22	0.08 – 0.12	[98]
2024	SiC	Fe <sub>70</sub> Ni <sub>30</sub>	SiV	7.5 – 7.9	1795 – 1805	8 – 12	198 – 202	[99]
				8 – 8.4	1995 – 2005	28 – 32	298 – 302	
				4.8 – 5.2	1414 – 1425	8 – 12	198 – 202	
				5 – 5.4	1495 – 1505	8 – 12	298 – 302	
2015	C <sub>10</sub> H <sub>8</sub>	Germanium	GeV	7 – 8	1295 – 1305	0.1 – 0.12	10 – 15	[64]
				8 – 9	1595 – 1605		0.03 – 0.1	
2019	C <sub>10</sub> H <sub>16</sub>	C <sub>24</sub> H <sub>20</sub> Ge	GeV	8.3 – 8.7	1895 – 1905	235 – 245	98 – 102	[79]
2021	C <sub>10</sub> H <sub>8</sub>	C <sub>24</sub> H <sub>20</sub> Ge	GeV	7.8 – 8.2	1295 – 1305	0.8 – 1.2	0.1 – 0.5	[101]
2021	Graphite and 0.5 mm Diamond	Germanium	GeV	6.8 – 7.2	1795 – 1805	3590 – 3605	448 – 452	[102]
2019	Graphite	Sn <sub>20</sub> Mg <sub>80</sub>	SnV	7.3 – 7.7	1415 – 1425	5910 – 4005	3995 – 4005	[101]
		Sn <sub>20</sub> Mg <sub>80</sub>			1895 – 1905	235 – 245	89 – 102	
		Sn <sub>50</sub> Mg <sub>50</sub>			1995 – 2005	475 – 485	118 – 122	
		Sn <sub>70</sub> Mg <sub>30</sub>			1845 – 1855	58 – 62	9 – 11	
		Sn <sub>20</sub> Mg <sub>80</sub>			1895 – 1905	235 – 245	98 – 102	
		Sn <sub>50</sub> Mg <sub>50</sub>			1795 – 1805	118 – 122	49 – 51	
		Sn <sub>70</sub> Mg <sub>30</sub>			1695 – 1705	475 – 485	2498 – 2502	
		Sn <sub>20</sub> Mg <sub>80</sub>			1895 – 1905	175 – 185	98 – 102	
		Sn <sub>50</sub> Mg <sub>50</sub>			1795 – 1805	118 – 122	49 – 51	
		Sn <sub>70</sub> Mg <sub>30</sub>			1695 – 1705	475 – 485	2498 – 2502	
		Sn-C system			1695 – 1705	475 – 485	2498 – 2502	
2021	Graphite	Sn-C system	SnV	9.2 – 9.6	1695 – 1705	475 – 485	2498 – 2502	[104]

**Table ST2.** High pressure high temperature (HPHT). Relevant information includes the carbon/color center sources and experimental conditions such as chamber pressure, temperature, and run time reported in the literature for NV, SiV, GeV and SnV. Pressure and temperature values are in the range 5–9.7 kPa and 895–1,995 °C, respectively. The size of the synthesized diamond NPs is in the range 0.010–6,005 μm. The HPHT database includes a total of 680 data points.

**Table ST3.** CVD/MPCVD, sub-database

Year	Carbon source	Dopant source	Substrate	Color centers	Pressure range (kPa)	Temperature range (°C)	Power (kW)	References
2019	CH <sub>4</sub> :H <sub>2</sub> (8:192 <i>sccm</i> )	N <sub>2</sub> (2 <i>sccm</i> )	CVD SCD on Si	NV	9.3 – 9.7	895 – 905	0.43 – 0.47	[55]
2020	CH <sub>4</sub> :H <sub>2</sub> (1:100 <i>sccm</i> )	N <sub>2</sub> (2–10 <i>sccm</i> )	CVD SCD on <i>p</i> -type Si	NV	3.9 – 4.1	745 – 755	0.7 – 0.9	[106]
2022	CH <sub>4</sub> (20 <i>sccm</i> )	N <sub>2</sub> O (480 <i>sccm</i> )	Ib-type HPHT (100) SCD	NV	19 – 21	995 – 1005	2.9 – 3.1	[29]
2022	CH <sub>4</sub> :H <sub>2</sub> (6.8:180 <i>sccm</i> )	N <sub>2</sub> (2 <i>sccm</i> )	Carbon films with SCDNs on (100) Si	NV	9.4 – 10.2	895 – 905	1.8 – 2.2	[60]
2023	CH <sub>4</sub> :H <sub>2</sub> (4: 200 <i>sccm</i> )	N <sub>2</sub> (5 <i>sccm</i> )	Ib-type HPHT (100) SCD	NV	10 – 12	1085 – 1095	1.3 – 1.5	[78]
2024	CH <sub>4</sub> :H <sub>2</sub> (3: 300 <i>sccm</i> )	C <sub>20</sub> H <sub>19</sub> N	Tantalum	NV	2.6 – 2.7	745 – 755	1.2 – 1.4	[105]
2024	CH <sub>4</sub> (9 <i>sccm</i> )	N <sub>2</sub> (150 <i>sccm</i> )	(001) CVD diamond	NV	0.9 – 1.1	945 – 955	5.8 – 6.2	[44]
2017	CH <sub>4</sub> :H <sub>2</sub> (24: 375 <i>sccm</i> )	SiH <sub>4</sub> (1 <i>sccm</i> )	HPHT thinfilm	SiV	17 – 18	845 – 855	3.1 – 3.5	[109], [111]
2017	CH <sub>4</sub> :H <sub>2</sub> (20: 378 <i>sccm</i> )	SiH <sub>4</sub> (0.8–2.4 <i>sccm</i> )	(100) CVD SCD d or Ib-type MPCVD SCD	SiV	2.9 – 3.3	345 – 355	0.3 – 0.5	[56]
2019	CH <sub>4</sub> :H <sub>2</sub> (8:192 <i>sccm</i> )	Si substrate	CVD SCD on Si	SiV	9 – 10	895 – 905	0.4 – 0.5	[55]
2021	CH <sub>4</sub> :H <sub>2</sub> (4: 400 <i>sccm</i> )	SiH <sub>4</sub> (5–30 <i>sccm</i> )	Si doped diamond or SiC	SiV	2.8 – 3.2	655 – 665	5 – 7	[108]
2021	CH <sub>4</sub> :H <sub>2</sub> (20: 400 <i>sccm</i> )	Si(CH <sub>3</sub> ) <sub>4</sub> (1 – 30 <i>sccm</i> )	(100) Si	SiV	5.5 – 6.5	855 – 865	5.8 – 6.2	[110]
2022	CH <sub>4</sub> :H <sub>2</sub> (8: 192 <i>sccm</i> )	Carbon films with SCDNs	(100) Si	SiV	9.6 – 10	655 – 665	1.8 – 2.2	[60]
2024	CH <sub>4</sub> :H <sub>2</sub> (24: 364 <i>sccm</i> )	SiH <sub>4</sub> (0.07 <i>sccm</i> )	Ila-type HPHT SCD	SiV	22 – 23	1045 – 1055	3.4 – 3.8	[112]
2024	CH <sub>4</sub> :H <sub>2</sub> (18: 300 <i>sccm</i> )	SiH <sub>4</sub> (1 <i>sccm</i> )	Ib-type MPCVD SCD	SiV	3.3 – 3.9	655 – 665	5 – 7	[113]
2015	CH <sub>4</sub> :H <sub>2</sub> (2:198 <i>sccm</i> )	Ge	Ib-type diamond, Element six	GeV	5 – 7	815 – 825	2 – 3	[70]
2018	CH <sub>4</sub> :H <sub>2</sub> (1: 100 <i>sccm</i> )	Ge/GeO <sub>2</sub>	(100) CVD SCD	GeV	7.8 – 8.2	345 – 355	0.85 – 0.95	[118]
2019	CH <sub>4</sub> :H <sub>2</sub> (8:192 <i>sccm</i> )	Ge	CVD SCD on Si	GeV	9.1 – 9.9	895 – 905	0.4 – 0.5	[55]
2019	CH <sub>4</sub> :H <sub>2</sub> (4: 400 <i>sccm</i> )	Ge/GeO <sub>2</sub>	SCD on sapphire	GeV	7.5 – 8.5	995 – 1005	0.8 – 1	[119]
2020	CH <sub>4</sub> :H <sub>2</sub> (1: 100 <i>sccm</i> )	GeO <sub>2</sub>	Sapphire	GeV	10 – 11	995 – 1005	0.8 – 1	[116]
2021	CH <sub>4</sub> :H <sub>2</sub> (20: 478 <i>sccm</i> )	GeH <sub>4</sub> (2 <i>sccm</i> )	HPHT diamond	GeV	11 – 11.6 17 – 18	845 – 855 985 – 995	2.4 – 2.8	[120]
2022	CH <sub>4</sub> :H <sub>2</sub> (8:192 <i>sccm</i> )	Ge	Carbon films with SCDNs on (100) Si	GeV	9.6 – 10	895 – 905	1.8 – 2.2	[60]
2022	CH <sub>4</sub> :H <sub>2</sub> (15:233 <i>sccm</i> )	GeH <sub>4</sub> (2 <i>sccm</i> )	(100) Si	GeV	14 – 15	795 – 805	4.7 – 5.1	[57]
2023	CH <sub>4</sub> :H <sub>2</sub> (20:480 <i>sccm</i> )	GeH <sub>4</sub> (0.3 <i>sccm</i> )	Ila-type diamond Element six-SCD	GeV	8 – 9	845 – 855	4.8 – 5.2	[117]
2020	CH <sub>4</sub> :H <sub>2</sub> (1: 100 <i>sccm</i> )	SnO <sub>2</sub> /SnCl <sub>2</sub>	Sapphire	SnV	10 – 11	995 – 1005	0.8 – 1	[116]
2020	CH <sub>4</sub> :H <sub>2</sub> (0.5: 300 <i>sccm</i> )	<sup>120</sup> Sn <sup>+</sup>	Electronic-grade diamond, Element Six	SnV	2.8 – 3.4	645 – 655	0.9 – 1.3	[62]

**Table ST3.** chemical vapor deposition (CVD) and Microwave Plasma Chemical Vapor Deposition (MCVD). Relevant information includes the carbon/color center sources, the type of substrate and experimental conditions such as chamber pressure, substrate temperature, and power of the microwave (MW) plasma reported in the literature for four different DVCCs: NV, SiV, GeV and SnV. Chamber pressure and substrate temperature values are in the range 0.9–18 kPa and 345–1,095 °C, respectively, while the power of the MW plasma is in the range 0.3–7 kW. The CVD/MPCVD-relevant database includes a total of 600 data points.

**Table ST4.** Ion implantation, sub-database

Year	Ion source	Substrate	Color centers	Implantation energy range(keV)	Fluences range (ions/cm <sup>2</sup> )	Temperature range (°C)	Annealing time range (min)	References	
2016	He <sup>+</sup>	Ila-type (001) MPCVD SCD	NV	3.9 – 4.1	9.0×10 <sup>10</sup> – 1.1×10 <sup>11</sup>	895 – 905	118 – 122	[146]	
2019	<sup>14</sup> N <sup>+</sup>	Ila- type (111) HPHT diamond	NV	29 – 31	9.0×10 <sup>10</sup> – 1.1×10 <sup>11</sup>	745 – 755	28 – 32	[28]	
2019	<sup>12</sup> C <sup>+</sup>	Ib-type SCD	NV	2.9 – 3.1	4.9×10 <sup>12</sup> – 5.1×10 <sup>12</sup>	1195 – 1205	118 – 122	[146]	
2019	N <sup>+</sup>	Quantum grade diamond, Element Six	NV	39 – 41	9.0×10 <sup>9</sup> – 1.1×10 <sup>10</sup>	795 – 805	239 – 241	[122]	
2020	N <sup>+</sup>	(001) CVD diamond	NV	34 – 36	2.9×10 <sup>9</sup> – 3.1×10 <sup>9</sup>	895 – 905	239 - 241	[123]	
2020	<sup>15</sup> N <sup>+</sup>	Electronic-grade (100) diamond-Element Six	NV	4 – 6	9.0×10 <sup>11</sup> – 1.1×10 <sup>12</sup>	995 – 1005	178 – 182	[124]	
2021	<sup>14</sup> N <sup>+</sup>	Ila-type (100) CVD diamond	NV	29 – 31	9.0×10 <sup>9</sup> – 1.1×10 <sup>10</sup>	785 – 805	118 – 122	[125]	
2022	<sup>12</sup> C <sup>+</sup>	Electronic grade diamond, Element-Six	NV	11 – 13	9.0×10 <sup>10</sup> – 1.1×10 <sup>10</sup>	785 – 805	118 – 122	[76]	
			NV	49 – 51	4.9×10 <sup>8</sup> – 5.1×10 <sup>8</sup>	795 – 805	118 – 122		
2023	<sup>15</sup> N <sup>+</sup>	Electronic grade diamond-Element Six	NV	5 – 7	9.0×10 <sup>10</sup> – 1.1×10 <sup>11</sup>	1195 – 1205	118 – 122	[126]	
2024	N <sub>2</sub> <sup>+</sup>	Ib-type MPCVD and HPHT diamond	NV	39 – 41	1.9×10 <sup>14</sup> – 2.1×10 <sup>14</sup>	795 – 805	118 – 122	[127]	
2016	<sup>29</sup> Si <sup>+</sup>	Ib-type CVD diamond Element-Six	SiV	145 – 155	9.0×10 <sup>9</sup> – 1.1×10 <sup>10</sup>	1195 – 1205	115 – 125	[128]	
2018	Si <sup>+</sup>	CVD SCD, Element-Six	SiV	395 – 405	6.5×10 <sup>10</sup> – 7.0×10 <sup>10</sup>	1195 – 1205	115 – 125	[129]	
2018	Si <sup>+</sup>	CVD SCD, Element-Six	SiV	70 – 74	3.0×10 <sup>13</sup> – 3.4×10 <sup>13</sup>	1095 – 1105	59 – 61	[130]	
2020	Si <sup>+</sup> /Si <sup>2+</sup> / <sup>3</sup> Si <sup>+</sup>	P-doped MPCVD diamond on HPHT	SiV	4.5 – 5.5	9.0×10 <sup>9</sup> – 1.1×10 <sup>10</sup>	1495 – 1505	85 – 95	[77]	
	<sup>28</sup> Si <sup>+</sup> / <sup>29</sup> Si <sup>+</sup> / <sup>30</sup> Si <sup>+</sup>	Ila-type electronic grade diamond, Element-Six		1.9×10 <sup>10</sup> – 2.1×10 <sup>10</sup>	1495 – 1505	85 – 95			
				4.9×10 <sup>11</sup> – 5.1×10 <sup>11</sup>					
				4.9×10 <sup>12</sup> – 5.1×10 <sup>12</sup>					
2018	Si <sup>3+</sup>	CVD SCD	SiV	9 – 11	9.0×10 <sup>11</sup> – 1.1×10 <sup>12</sup>	1045 – 1055	115 – 125	[72]	
					4.9×10 <sup>12</sup> – 5.1×10 <sup>12</sup>				
					9.0×10 <sup>12</sup> – 1.1×10 <sup>13</sup>				
					9.0×10 <sup>13</sup> – 1.1×10 <sup>14</sup>				
				19 – 21	9.0×10 <sup>11</sup> – 1.1×10 <sup>12</sup>	1045 – 1055			
					4.9×10 <sup>12</sup> – 5.1×10 <sup>12</sup>				
					9.0×10 <sup>12</sup> – 1.1×10 <sup>13</sup>				
					9.0×10 <sup>13</sup> – 1.1×10 <sup>14</sup>				
				49 – 51	9.0×10 <sup>11</sup> – 1.1×10 <sup>12</sup>	1045 – 1055			
					4.9×10 <sup>12</sup> – 5.1×10 <sup>12</sup>				
					9.0×10 <sup>12</sup> – 1.1×10 <sup>13</sup>				
					9.0×10 <sup>13</sup> – 1.1×10 <sup>14</sup>				
				99 – 101	9.0×10 <sup>11</sup> – 1.1×10 <sup>12</sup>	1045 – 1055			
					4.9×10 <sup>12</sup> – 5.1×10 <sup>12</sup>				
					9.0×10 <sup>12</sup> – 1.1×10 <sup>13</sup>				
					9.0×10 <sup>13</sup> – 1.1×10 <sup>14</sup>				
2023	<sup>29</sup> Si <sup>+</sup>	CVD SCD, Element-Six	SiV	65 – 75	9.0×10 <sup>8</sup> – 1.1×10 <sup>9</sup>	1195 – 1205	115 – 125	[132]	
2024	<sup>28</sup> Si <sup>+</sup>	HPHT and CVD SCD	SiV	17995 – 18005	6.0×10 <sup>12</sup> – 7.0×10 <sup>12</sup>	1795 – 1805	0.9 – 1.1	[133]	
2024	<sup>28</sup> Si <sup>+</sup>	CVD SCD, Element-Six	SiV	130 – 134	2.9×10 <sup>11</sup> – 3.1×10 <sup>11</sup>	1095 – 1105	595 – 605	[134]	
2018	Si <sup>+</sup>	CVD SCD, Element Six	SiV	95 – 97	3.0×10 <sup>13</sup> – 3.4×10 <sup>13</sup>	1095 – 1105	118 – 122	[130]	
				2995 – 3005	6.5×10 <sup>10</sup> – 7.0×10 <sup>10</sup>		115 – 125		
2018	Si <sup>3+</sup>	CVD SCD	SiV	19 – 21	9.0×10 <sup>7</sup> – 1.1×10 <sup>8</sup>	1195 – 1205	59 – 61		
2020	Si <sup>+</sup> /Si <sup>2+</sup> / <sup>3</sup> Si <sup>+</sup>	P-doped MPCVD diamond on HPHT	SiV	19 – 21	9.0×10 <sup>11</sup> – 1.1×10 <sup>12</sup>	1195 – 1205	59 – 61	[131]	
					9.0×10 <sup>12</sup> – 1.1×10 <sup>13</sup>				
					9.0×10 <sup>13</sup> – 1.1×10 <sup>14</sup>				
2015	Ge <sup>+</sup>	Ila-type (001) SCD Element-Six	GeV	145 – 155	3.0×10 <sup>8</sup> – 4.0×10 <sup>8</sup>	795 – 805	25 – 35	[70]	
					5.0×10 <sup>13</sup> – 6.5×10 <sup>13</sup>				
				250 – 270	3.0×10 <sup>8</sup> – 4.0×10 <sup>8</sup>				
					5.0×10 <sup>13</sup> – 6.5×10 <sup>13</sup>				
2021	Ge <sup>+</sup>	MPCVD SCD on Si	GeV	95 – 105	9.0×10 <sup>13</sup> – 1.1×10 <sup>14</sup>	895 – 905	115 - 125	[135]	
2024	<sup>75</sup> Ge <sup>+</sup>	(100) SCD-Element-Six	GeV	25 – 35	1.5×10 <sup>12</sup> – 2.5×10 <sup>12</sup>	895 – 905	8 – 12	[75]	
					1.5×10 <sup>13</sup> – 2.5×10 <sup>13</sup>				
2024	<sup>74</sup> Ge <sup>+</sup>	HPHT and CVD SCD	GeV	17000 – 19000	1.4×10 <sup>13</sup> – 1.6×10 <sup>13</sup>	595 – 605	595 – 605	[133]	
						695 – 705	595 – 605		
						1700 – 1900	0.5 – 1.5		
2017	Sn <sup>+</sup>	CVD SCD	SnV	128 – 132	1.9×10 <sup>8</sup> – 2.1×10 <sup>8</sup>	795 – 805	28 – 32	[136]	
				148 – 152	1.9×10 <sup>13</sup> – 2.1×10 <sup>13</sup>	2095 – 2105			
2017	Sn <sup>+</sup>	Ila-type SCD Element-Six	SnV	58 – 62	2.9×10 <sup>11</sup> – 3.1×10 <sup>11</sup>	895 – 905	28 – 32	[73]	
					1.9×10 <sup>10</sup> – 2.1×10 <sup>10</sup>	945 – 955	115 – 125		
					4.9×10 <sup>10</sup> – 5.1×10 <sup>10</sup>				
					9.0×10 <sup>10</sup> – 1.1×10 <sup>11</sup>				
				9995 – 10005	4.9×10 <sup>13</sup> – 5.1×10 <sup>13</sup>				
2019	Sn <sup>+</sup>	Quantum grade diamond-Element Six	SnV	78 – 82	2.9×10 <sup>11</sup> – 3.1×10 <sup>11</sup>	1195 – 1205	475 – 485	[122]	
2019	<sup>120</sup> Sn <sup>+</sup>	CVD SCD-Element Six	SnV	365 – 375	1.9×10 <sup>11</sup> – 2.1×10 <sup>11</sup>	1095 – 1105	88 – 92	[87]	
2020	<sup>120</sup> Sn <sup>+</sup>	Electronic-grade diamond Element Six	SnV	365 – 375	1.9×10 <sup>13</sup> – 2.1×10 <sup>13</sup>	1095 – 1105	28 – 32	[62]	
2020	Sn <sup>+</sup>	(001) CVD SCD	SnV	75 – 77	9.0×10 <sup>9</sup> – 1.1×10 <sup>10</sup>	1195 – 1205	235 – 245	[123]	
					9.0×10 <sup>10</sup> – 1.1×10 <sup>11</sup>				
2020	<sup>121</sup> Sn <sup>+</sup>	Ila-type (110) Dutch Diamond Technologies	SnV	58 – 62	1.9×10 <sup>12</sup> – 2.1×10 <sup>12</sup>	915 – 925	8 – 12	[138]	
2020	Sn <sup>+</sup>	Electronic-grade (001) diamond	SnV	79 – 81	1.9×10 <sup>13</sup> – 2.1×10 <sup>13</sup>	1195 – 1205	235 – 245	[137]	
	<sup>120</sup> Sn <sup>+</sup>	Electronic-grade diamond-Element Six		695 – 705		2095 – 2105	28 – 32		
2021	Sn <sup>+</sup>	Ila-type (100) diamond	SnV	58 – 62	1.9×10 <sup>10</sup> – 2.1×10 <sup>10</sup>	1195 – 1205	115 – 125	[121]	
2021	Sn <sup>+</sup>	Ila-type CVD SCD-Element Six	SnV	345 – 355	9.0×10 <sup>8</sup> – 1.1×10 <sup>9</sup>	1195 -1205	118 – 122	[65]	
2022	Sn <sup>+</sup>	Electronic-grade (001) diamond Element-Six	SnV	695 – 705	9.0×10 <sup>8</sup> – 1.1×10 <sup>9</sup>	2095 – 2105	115 – 125	[90]	
					9.0×10 <sup>11</sup> – 1.1×10 <sup>12</sup>				
					7.9×10 <sup>13</sup> – 8.1×10 <sup>13</sup>				
2024	<sup>117</sup> Sn <sup>+</sup>	Electronic-grade diamond-Element Six	SnV	345 - 355	9.0×10 <sup>10</sup> – 1.1×10 <sup>11</sup>	1195 – 1205	715 – 725	[74]	

**Table ST4.** Ion implantation. Relevant information includes the ion source, substrate, implantation energy, fluence, temperature and annealing time reported for the different DVCCs. Implantation energies and fluences are in the range 3.9–19,000 keV and 1.9×10<sup>8</sup>–2.1×10<sup>14</sup> ions/cm<sup>2</sup>, respectively; the temperature and annealing time are in the range 595–2,105 °C and 0.9–725 minutes, respectively. The Ion irradiation-relevant database includes a total of 284 data points.

**Table ST5.** Electron/ion irradiation, sub-database

Year	Ion source	Substrate	Color centers	Implantation energy range (MeV)	Fluences range (ions/cm <sup>2</sup> )	Temperature range (°C)	Annealing time range (min)	References
2017	Electron beam	CVD SCD, Element-Six	NV	4.3 – 4.7	5.0×10 <sup>18</sup> – 6.0×10 <sup>18</sup>	1195 – 1205	118 – 122	[85]
2019	Electron beam	CVD SCD	NV	1.8 – 2.2	4.9×10 <sup>15</sup> – 5.1×10 <sup>13</sup>	1095 – 1105	598 – 602	[141]
2020	Electron beam	SCD-Sumitomo Electric	NV	4.4 – 4.8	9.0×10 <sup>17</sup> – 1.1×10 <sup>18</sup>	795 – 805	58 – 62	[140]
2021	Au <sup>+</sup> ion	Ib-type diamond-Element Six	NV	2.2 – 2.6	7.9×10 <sup>16</sup> – 8.1×10 <sup>16</sup>	1045 – 1055	148 – 152	[95]
					3.9×10 <sup>17</sup> – 4.1×10 <sup>17</sup>			
					3.9×10 <sup>18</sup> – 4.1×10 <sup>18</sup>			
					9.0×10 <sup>18</sup> – 1.0×10 <sup>19</sup>			
2022	Electron beam	MPCVD SCD on Ib-type HPHT (001) diamond	NV	1.9 – 2.1	9.0×10 <sup>15</sup> – 1.1×10 <sup>16</sup>	995 – 1005	118 - 122	[78]
					1.9×10 <sup>16</sup> – 2.1×10 <sup>16</sup>			
					9.0×10 <sup>16</sup> – 1.1×10 <sup>17</sup>			
					1.9×10 <sup>17</sup> – 2.1×10 <sup>17</sup>			
2022	Electron beam	N <sub>2</sub> MPCVD SCD on Ila-type HPHT diamond	NV	0.9 – 1.1	9.0×10 <sup>16</sup> – 1.1×10 <sup>17</sup>	995 – 1005	118 – 122	[41]
					2.9×10 <sup>17</sup> – 3.1×10 <sup>17</sup>			
					9.0×10 <sup>17</sup> – 1.1×10 <sup>18</sup>			
					2.9×10 <sup>18</sup> – 3.1×10 <sup>18</sup>			
2022	Electron beam	HPHT (100) diamond	NV	1.9 – 2.1	9.0×10 <sup>17</sup> – 1.1×10 <sup>18</sup>	995 – 1005	118 – 122	[42]
				1.8 – 2.2	9.0×10 <sup>16</sup> – 1.1×10 <sup>17</sup>			
2022	Electron beam	Ib-type (100) SCD	NV	2.1 – 2.5	2.9×10 <sup>18</sup> – 3.1×10 <sup>18</sup>	795 – 805	118 – 122	[29]
				2.0 – 2.6		795 – 805	118 – 122	
						98 – 102	58 – 62	
						1195 – 1205	58 – 62	
2023	Electron beam	N <sub>2</sub> MPCVD SCD	NV	9 – 11	2.2×10 <sup>18</sup> – 2.6×10 <sup>18</sup>	845 – 855	238 – 242	[30]
2023	Electron beam	HPHT SCD-Element Six	NV	0.18 – 0.22	1.4×10 <sup>18</sup> – 1.6×10 <sup>18</sup>	795 – 805	238 – 242	[142]
					1.9×10 <sup>18</sup> – 2.1×10 <sup>18</sup>			
					9.0×10 <sup>18</sup> – 1.1×10 <sup>19</sup>			
					4.9×10 <sup>19</sup> – 5.1×10 <sup>19</sup>			
					9.0×10 <sup>19</sup> – 1.1×10 <sup>20</sup>			
					1.9×10 <sup>20</sup> – 2.1×10 <sup>20</sup>			
					4.9×10 <sup>20</sup> – 5.1×10 <sup>20</sup>			
				153 – 157	1.4×10 <sup>18</sup> – 1.6×10 <sup>18</sup>			
2024	Electron beam	Fe-Co-Ni HPHT diamond	NV	9 – 11	9.0×10 <sup>16</sup> – 1.1×10 <sup>17</sup>	795 – 805	58 – 62	[31]
2021	Electron beam	5nm DNDs	SiV	1.8 – 2.2	4.9×10 <sup>18</sup> – 5.1×10 <sup>18</sup>	1095 – 1105	59 – 61	[144]
2023	Electron beam	Ila-type HPHT SCD	SiV	1.8 – 2.2	1.9×10 <sup>18</sup> – 2.1×10 <sup>18</sup>	1595 – 1605	59 – 61	[86]
2024	He <sup>+</sup> ion	MPCVD SCD	SiV	1.5 – 1.9	2.0×10 <sup>15</sup> – 3.0×10 <sup>15</sup>	795 – 805	118 – 122	[112]
2021	Ge <sup>+</sup> ion	MPCVD SCD on (100) Si	GeV	0.08 – 0.12	9.0×10 <sup>14</sup> – 1.1×10 <sup>15</sup>	895 – 905	118 – 122	[135]

**Table ST5.** Electron irradiation. Relevant information includes the ion source, substrate, implantation energy and fluence, as well as temperature and annealing time reported in the literature for three different DVCCs: NV, SiV and GeV. Implantation energies and fluences are in the range 0.08–157 MeV and 4.9×10<sup>13</sup>–5.1×10<sup>20</sup> ions/cm<sup>2</sup>, respectively; the temperature and annealing time are in the range 98–1,605 °C and 58–602 minutes, respectively. The Ion irradiation-relevant database includes a total of 128 data points.

### S3. Hyperparameters of the DTR and XGB models

Table ST6 lists the hyperparameters with their value ranges, for the DTR and XGB models, including the subsample ratio instance, tree depth, sampling, and learning rate of full training. The hyperparameters alpha (L1) and lambda (L2) are regularization terms that appear in the regularization function  $\Omega(f_t)$  in eq. (4) of the main text and are applied to control the complexity of the model and to avoid overfitting.

**Table ST6.** Hyperparameters for the DTR/XGB predictors

Model	Hyperparameter	Search range	Optimum value/feature	Model	Hyperparameter	Search range	Optimum value/feature
Decision-Tree Regressor	Criterion	mse, friedman_mse, mae	friedman_mse	XGBBoosting	booster	gbtree, gblinear	gbtree
	splitter	best, random	random		objective	reg, rank, binary, survival, multi ...	reg: linear
	max_features	sqrt, auto, None	None		learning_rate	0 – 1	0.06
	max_depth	10 – 100	10		Eval_metric	logloss	rmse
	min_samples_split	2 – 10	2		max_depth	-	1
	min_samples_leaf	1 – 5	1		subsample	0 – 1	0.8
	min_weight_fraction_leaf	0 – 1000	0		Colsample_bytree	0 – 1	1
	max_leaf_nodes	None, random	None		lambda	0 – 1000	10
	random_state	None, random	42		alpha	0 – 1000	0
					seed	0 – 100	66

**Table ST6.** List and value ranges of the hyperparameters of the trained DTR and XGB predictors.

## S4. Performance of the DTR and XGB Models

Tables ST7 and ST8 below summarize the values of the statistical indicators,  $R^2$ , MSE and MAE we used to assess the performance of the DTR and XGB predictors for *Task I* and *II*, for each technique: HPHT, CVD/MPCVD, Ion implantation and Electron/ion irradiation, respectively (see main text).

**Table ST7. *Task I*:** DTR/XGB performance for each synthesis method

		$R^2$		MSE		MAE	
	DVCC	DTR	XGB	DTR	XGB	DTR	XGB
HPHT	NV	(1.000) 1.000	(1.000) 1.000	(0.0044) 0.0093	(0.0000) 0.0043	(0.0449) 0.0750	(0.0007) 0.0582
	SiV	(1.000) 1.000	(1.000) 1.000	(0.1964) 0.2174	(0.0003) 0.1232	(0.2327) 0.2625	(0.0084) 0.0865
	GeV	(1.000) 1.000	(1.000) 1.000	(0.0186) 0.0112	(0.0000) 0.0023	(0.0751) 0.0679	(0.0023) 0.0330
	SnV	(1.000) 1.000	(1.000) 1.000	(0.0063) 0.0392	(0.0001) 0.0287	(0.0528) 0.1659	(0.0068) 0.1349
CVD/MPCVD	NV	(0.989) 0.943	(1.000) 0.986	(0.0295) 0.2249	(0.0011) 0.0555	(0.0406) 0.1739	(0.0198) 0.0814
	SiV	(1.000) 1.000	(1.000) 0.994	(0.0014) 0.0010	(0.0024) 0.0287	(0.0158) 0.0239	(0.0263) 0.0860
	GeV	(0.989) 0.975	(1.000) 0.966	(0.0270) 0.0587	(0.0008) 0.0792	(0.0455) 0.0693	(0.0135) 0.0846
	SnV	(1.000) 0.998	(1.000) 0.999	(0.0000) 0.0000	(0.0000) 0.0000	(0.0010) 0.0048	(0.0009) 0.0038
Ion implantation	NV	(1.000) 1.000	(1.000) 1.000	(0.0000) 0.0652	(0.0000) 0.1230	(0.0000) 0.2202	(0.0017) 0.2708
	SiV	(1.000) 1.000	(1.000) 0.999	(0.0000) 0.5179	(0.0000) 4.9399	(0.0000) 0.3599	(0.0019) 0.5882
	GeV	(1.000) 1.000	(1.000) 1.000	(0.0000) 0.0815	(0.0000) 0.1613	(0.0000) 0.1658	(0.0015) 0.3174
	SnV	(1.000) 1.000	(1.000) 1.000	(0.0000) 0.0741	(0.0000) 0.2941	(0.0000) 0.2054	(0.0028) 0.3237
Electron/ion irradiation	NV	(1.000) 1.000	(1.000) 1.000	(0.0000) 0.2370	(0.0000) 0.3753	(0.0000) 0.4741	(0.0013) 0.5557
	SiV	(1.000) 1.000	(1.000) 1.000	(0.0000) 0.0148	(0.0000) 0.0132	(0.0000) 0.0370	(0.0003) 0.0302
	GeV	(1.000) 0.973	(1.000) 0.946	(0.0000) 0.0443	(0.0000) 0.0883	(0.0000) 0.2105	(0.0008) 0.2802

**Table ST7.** List of the statistical indicators,  $R^2$ , MSE and MAE for *Task I*, for the different material synthesis techniques (HPHT, CVD/MPCVD, Ion implantation and Electron/ion irradiation), organized by color centers (NV, SiV, GeV and SnV) and by method (DTR and XGB). Values for the test dataset are shown in red, while values for the training dataset are shown in black and within parentheses. The values for the training dataset are merely used for control.

**Table ST8.** *Task II:* DTR/XGB performance for each synthesis method

		<b>R<sup>2</sup></b>		<b>MSE</b>		<b>MAE</b>	
	<b>DVCC</b>	<b>DTR</b>	<b>XGB</b>	<b>DTR</b>	<b>XGB</b>	<b>DTR</b>	<b>XGB</b>
<b>HPHT</b>	<b>NV</b>	(1.000) 0.997	(1.000) 0.996	(0.0000) 0.0014	(0.0000) 0.0021	(0.0000) 0.0361	(0.0010) 0.0413
	<b>SiV</b>	(1.000) 1.000	(1.000) 0.942	(0.0000) 0.0001	(0.0000) 0.0212	(0.0000) 0.0058	(0.0006) 0.0296
	<b>GeV</b>	(1.000) 0.998	(1.000) 0.706	(0.0000) 0.0011	(0.0000) 0.1519	(0.0000) 0.0311	(0.0006) 0.1962
	<b>SnV</b>	(1.000) 1.000	(1.000) 0.987	(0.0000) 0.0002	(0.0000) 0.0094	(0.0000) 0.0147	(0.0005) 0.0483
<b>CVD/MPCVD</b>	<b>NV</b>	(1.000) 0.993	(1.000) 0.841	(0.0000) 0.0028	(0.0000) 0.0685	(0.0000) 0.0507	(0.0007) 0.1747
	<b>SiV</b>	(1.000) 0.949	(1.000) 0.705	(0.0000) 0.0260	(0.0000) 0.1509	(0.0000) 0.0394	(0.0006) 0.1989
	<b>GeV</b>	(1.000) 1.000	(1.000) 0.988	(0.0000) 0.0001	(0.0000) 0.0049	(0.0000) 0.0098	(0.0006) 0.0295
	<b>SnV</b>	(1.000) 0.908	(1.000) 0.898	(0.0000) 0.0192	(0.0000) 0.0214	(0.0000) 0.1275	(0.0005) 0.1372
<b>Ion implantation</b>	<b>NV</b>	(1.000) 0.971	(1.000) 0.853	(0.0000) 0.0094	(0.0000) 0.0479	(0.0000) 0.0780	(0.0007) 0.1351
	<b>SiV</b>	(1.000) 0.941	(1.000) 0.905	(0.0000) 0.0257	(0.0000) 0.0413	(0.0000) 0.0423	(0.0011) 0.1267
	<b>GeV</b>	(1.000) 1.000	(1.000) 0.855	(0.0000) 0.0001	(0.0000) 0.0658	(0.0000) 0.0098	(0.0006) 0.1123
	<b>SnV</b>	(1.000) 0.998	(1.000) 0.956	(0.0000) 0.0006	(0.0000) 0.0104	(0.0000) 0.0223	(0.0009) 0.0670
<b>Electron/ion irradiation</b>	<b>NV</b>	(1.000) 0.949	(1.000) 0.911	(0.0000) 0.0162	(0.0000) 0.0284	(0.0000) 0.0388	(0.0017) 0.0833
	<b>SiV</b>	(1.000) 1.000	(1.000) 0.943	(0.0000) 0.0000	(0.0000) 0.0215	(0.0000) 0.0039	(0.0008) 0.0286
	<b>GeV</b>	(1.000) 0.987	(1.000) 0.981	(0.0000) 0.0045	(0.0000) 0.0065	(0.0000) 0.0632	(0.0007) 0.0720

**Table ST8.** List of the statistical indicators, R<sup>2</sup>, MSE and MAE for *Task II* for the different material synthesis techniques (HPHT, CVD/MPCVD, Ion implantation and Electron/ion irradiation), organized by color centers (NV, SiV, GeV and SnV) and by method (DTR and XGB). Values for the test dataset are shown in red, while values for the training dataset are shown in black and within parentheses. The values for the training dataset are merely used for control.

## S5. Ranges of DWF values

Table ST9 below shows the ranges of Debye-Waller factor (DWF) values measured for each color center, differentiated by synthesis technique.

**Table ST9.** Measured DWF values for each technique

	HPHT	CVD/MPCVD	Ion implantation	Electron/ion irradiation
NV	0.029 – 0.036	0.030 – 0.048	0.032 – 0.049	0.028 – 0.054
SiV	0.640 – 0.741	0.642 – 0.725	0.644 – 0.752	0.668 – 0.727
GeV	0.564 – 0.616	0.550 – 0.619	0.553 – 0.613	0.620 – 0.622
SnV	0.471 – 0.510	0.474 – 0.479	0.470 – 0.512	–

**Table ST9.** Ranges of measured DWF values organized by synthesis techniques, HPHT, CVD/MPCVD, Ion implantation and Electron/ion irradiation (columns) and organized by color centers, NV, SiV, GeV and SnV (rows). In red we highlight, for each color center, the techniques that yield the largest DWF value, which is usually desirable for optical, photonic and quantum applications.