




## RESEARCH ARTICLE

10.1029/2023GC011224

## Volcanic Ash Classification Through Machine Learning

Damià Benet<sup>1,2,3</sup> , Fidel Costa<sup>1</sup>, and Christina Widiwijayanti<sup>2</sup>

## Key Points:

- Volcanic ash particles are classified through machine learning algorithms into juvenile, lithic, free-crystal and altered material types
- Discriminant features per each particle type are revealed by the Shapley values of XGBoost's predictions
- Classification by a Vision Transformer model is very accurate and could be used by volcano observatories

## Supporting Information:

Supporting Information may be found in the online version of this article.

## Correspondence to:

D. Benet,  
dbenet@ipgp.fr

## Citation:

Benet, D., Costa, F., & Widiwijayanti, C. (2024). Volcanic ash classification through machine learning. *Geochemistry, Geophysics, Geosystems*, 25, e2023GC011224. <https://doi.org/10.1029/2023GC011224>

Received 6 SEP 2023

Accepted 22 DEC 2023

<sup>1</sup>Institut de Physique du Globe de Paris, Université Paris Cité, CNRS, Paris, France, <sup>2</sup>EOS, Earth Observatory of Singapore, Nanyang Technological University, Singapore, Singapore, <sup>3</sup>Asian School of the Environment, Nanyang Technological University, Singapore, Singapore

**Abstract** Volcanic ash provides information that can help understanding the evolution of volcanic activity during the early stages of a crisis and possible transitions toward different eruptive styles. Ash consists of particles from a range of origins within the volcanic system and its analysis can be indicative of the processes driving the eruptive activity. However, classifying ash particles into different types is not straightforward. Diagnostic observations for particle classification are not standardized and vary across samples. Here we explore the use of machine learning (ML) to improve the classification accuracy and reproducibility. We use a curated database of ash particles (VolcAshDB) to optimize and train two ML-based models: Extreme Gradient Boosting (XGBoost) that uses the measured physical attributes of the particles, from which predictions are interpreted by the SHapley Additive exPlanations (SHAP) method, and a Vision Transformer (ViT) that classifies binocular, multi-focused, particle images. We find that the XGBoost has an overall classification accuracy of 0.77 (*macro F1-score*), and specific features of color (*hue\_mean*) and texture (*correlation*) are the most discriminant between particle types. Classification using the particle images and the ViT is more accurate (*macro F1-score* of 0.93), with performances varying from 0.85 for samples of dome explosions, to 0.95 for phreatic and subplinian events. Notwithstanding the success of the classification algorithms, the training dataset is limited in number of particles, ranges of eruptive styles, and volcanoes. Thus, the algorithms should be tested further with additional samples, and it is likely that classification for a given volcano is more accurate than between volcanoes.

### 1. Introduction

A central challenge in volcanology is to anticipate the likely evolution of a restless volcano at a given point in time (Bebbington & Jenkins, 2019). During a period of unrest, small explosions or phreatic events may precede larger ones, or the volcano may remain at low activity levels and go back to dormancy (Marzocchi et al., 2012; Moran et al., 2011; Tilling, 2008). Moreover, many eruptions consist of various phases, changing or alternating between explosive and effusive eruptive styles over time. To evaluate whether a volcano will progress toward one type of activity or another, an array of geophysical and geochemical tools is used to monitor and interpret the processes happening underneath the volcano (Newhall & Punongbayan, 1996). However, interpretation may not be straightforward and available data are limited, and thus diagnosis is typically quite uncertain (Tilling, 2008).

An additional tool that can provide critical insights into the state of a volcano is studying the volcanic ash. Ash can be classified into particle types, also known as components, that are indicative of processes driving the activity (Alvarado et al., 2016; D'Oriano et al., 2022; Gaunt et al., 2016; Pardo et al., 2014; Re et al., 2021). For instance, juvenile particles are associated with the fragmentation of ascending magma at shallow depths, and their identification, together with other monitoring signals, may warn of an ensuing magmatic eruption. For example, a posteriori studies of ash from early and small phreatic eruptions of Mount St. Helens (USA), 1980 and Mount Unzen (Japan), 1991, identified minor number of juvenile particles in these pre-climactic deposits (Cashman & Hoblitt, 2004; Watanabe et al., 1999). Thus, had these been found in a timely manner, it could have altered the perception for explosive potential that followed (Cashman & Hoblitt, 2004). In other cases, the ambiguity of classification of the juvenile component in early explosions has led to very complex management of the volcanic crises such as the 1975–1977 Soufrière Guadeloupe crisis (Feuillard et al., 1983; Hincks et al., 2014; Le Guern et al., 1980). Furthermore, tracking the proportions of the different components in ash, their shape, and crystallinity, can give clues on possible transitions of eruption styles to better mitigate the associated hazards (e.g., Benet et al., 2021; Suzuki et al., 2013; Taddeucci et al., 2002).

© 2024 The Authors. *Geochemistry, Geophysics, Geosystems* published by Wiley Periodicals LLC on behalf of American Geophysical Union. This is an open access article under the terms of the [Creative Commons Attribution-NonCommercial-NoDerivs License](https://creativecommons.org/licenses/by/4.0/), which permits use and distribution in any medium, provided the original work is properly cited, the use is non-commercial and no modifications or adaptations are made.

The classification of particles into different types is typically done by collecting qualitative or quantitative data on a single particle level using a variety of techniques. This includes using a binocular microscope (e.g., D'Oriano et al., 2014; Miwa et al., 2009; Pardo et al., 2014) to observe the gloss, color and shape, as well as the particles' surface and shape (Dellino & La Volpe, 1996; Dürig et al., 2021; E. J. Liu et al., 2015; Ross et al., 2022). More detailed observations including the internal microstructures are typically done using the Scanning Electron Microscope (e.g., Miwa et al., 2013; Pardo et al., 2020), whereas the chemical analyses are conducted with the electron microprobe (Pardo et al., 2014), mass spectrometers (Rowe et al., 2008), and measurement of refractive indices (e.g., by the thermal immersion method; Watanabe et al., 1999). However, systematic and reproducible particle classification is challenging because there is no widespread consensus on the particle features that are diagnostic. These features may vary from sample to sample, depending on the eruptive style and the volcano (e.g., Pardo et al., 2014). Although a useful methodology to study juvenile pyroclasts has been proposed by Ross et al. (2022), the specific diagnostic features that separate different particle types are still unclear and difficult to apply consistently across a wide range of samples.

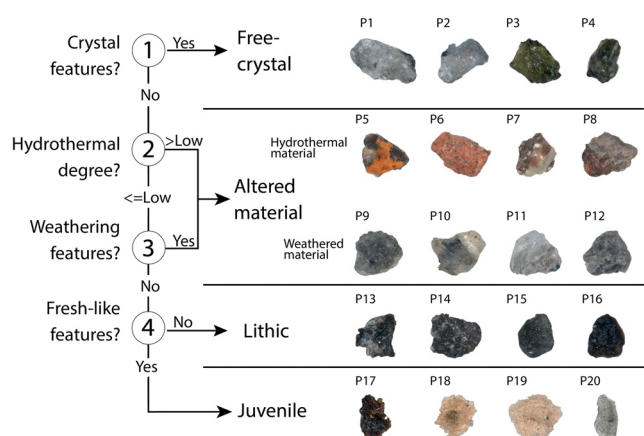
An approach commonly employed to address such classification challenges in various domains is through the utilization of Machine Learning (ML). ML-based models can classify complex images in a wide range of situations (He et al., 2015). ML-based models are capable of learning patterns to classify objects and use them for classification of future datasets, such as mushrooms (Lee et al., 2022) or leaf diseases (Sujatha et al., 2021). The study most closely related to ours is Shoji et al. (2018), where the authors successfully employed a neural network to classify volcanic ash particle shapes from samples from various eruptive activity types. In this work, we trained two models using the VolcAshDB curated dataset (Benet et al., 2024) with the objectives of (a) identifying the most important features for the discrimination of particle types and (b) obtaining a particle classifier as accurate as possible. The results of this study should be a step forward toward a universal and unbiased classification of ash particles as more data becomes available and better algorithms are developed.

## 2. Materials and Methods

### 2.1. VolcAshDB Data Set

We used the data from the open-access database VolcAshDB, which comprises images and measurements (here referred as features) of more than 6,300 volcanic ash particles (<http://volcashdb.ipgp.fr/>). These were obtained using the binocular microscope and processed to obtain multi-focused, high-resolution images (Benet et al., 2024). The images have been classified with a dichotomous key (Figure 1), using some key observational features as reported in Benet et al. (2024) and summarized in Table S1 in Supporting Information S1. The database contains ash particles from 12 samples from 8 volcanoes and 11 eruptions from a range of magma compositions and eruptive activity types (Table 1). These include (a) phreatic eruptions of Soufrière de Guadeloupe (Lesser Antilles) in 1976 and 1977 (Feuillard et al., 1983) and the early activity of April 1991 of Mt. Pinatubo (Philippines; Paladio-Melasantos et al., 1996), and Ontake (Japan) in 2014 (Miyagi et al., 2020), (b) dome explosions of Nevados de Chillán volcanic complex (Chile) from the beginning of the eruptive period in December 2016 and after the extrusion of a dome in April 2018 (Benet et al., 2021), explosions from Merapi volcano (Indonesia) in July and November 2013 (Nurfiani & Bouvet de Maisonneuve, 2018), (c) the basaltic lava fountaining of Cumbre Vieja (Canary Islands) in October 2021 (Romero et al., 2022), and (d) two samples from different locations (KE-DB2 and KE-DB3) of the plinian/sub-plinian eruptions of Kelud (Indonesia) in 2014 (Maeno et al., 2019; Utami et al., 2021), and a sample from the climactic plinian eruption of Mount St. Helens (USA) in 1980 (Scheidegger et al., 1982).

In addition to ash images, VolcAshDB includes (a) the value of 33 features of each ash particle related to its shape, texture, and color, (b) a label with the identification of the types of particles (free-crystal, altered material, juvenile, and lithic; Figure 1), and (c) metadata for each particle, such as the sample grain-size fraction, the number of magnifications used for image acquisition, amongst others. The shape features in the database have been used in previous studies (Cioni et al., 2014; Dellino & La Volpe, 1996; Dürig et al., 2018; Leibrandt & Le Pennec, 2015; E. J. Liu et al., 2015), and include those sensitive to particle-scale cavities (e.g., *solidity*), perimeter-based irregularities (e.g., *convexity*), and form (e.g., *elongation*; Liu et al., 2015). The textural features in VolcAshDB were obtained from calculations of the distribution of pixel intensities in grayscale across several particle regions



**Figure 1.** Example of classification flow and particle images in VolcAshDB based on the steps for petrographic classification in Benet et al. (2024). Note that the particle type altered material comprises both hydrothermal and weathered materials. The definition of the key diagnostic features used at each step can be found in Table S1 in Supporting Information S1.

based on the so-called Gray-Level Cooccurrence Matrix (GLCM, Haralick et al., 1973). From the GLCMs, we obtained features that indicate a more uniform texture (e.g., *homogeneity*), and those that indicate a more complex or heterogeneous texture (e.g., *dissimilarity*; Hall-Beyer, 2017). The color features of each particle were taken from the measurement of the mean, mode and standard deviation of the histogram distribution for each of the six channels in the Red-Green-Blue (RGB), and Hue-Saturation-Value (HSV) color spaces. For more details on the calculation and references of each feature, the reader is referred to Benet et al. (2024), and they are summarized with the abbreviation in Tables S2 and S3 in Supporting Information S1.

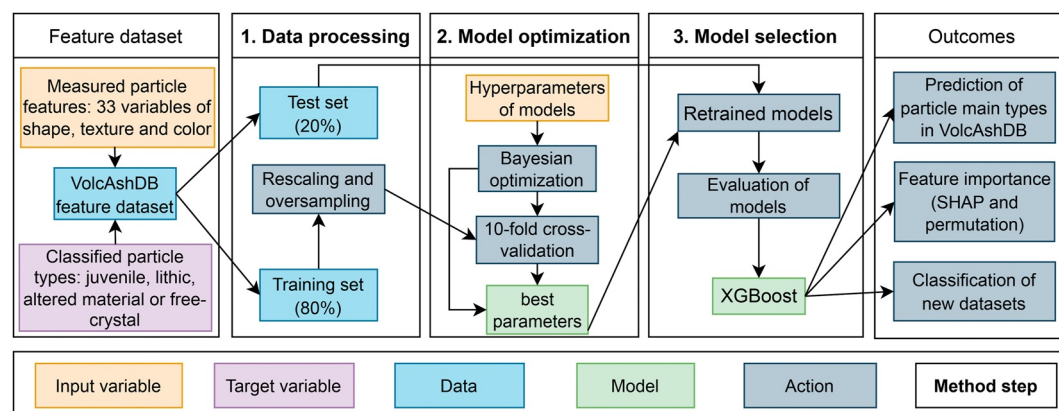
## 2.2. Development of a Particle Classifier Using the Measured Particle Features

The steps needed to develop a volcanic ash particle classifier vary if the input data are the measured features or the particle images directly. Because the particle types are already classified, the models are trained by supervised learning (Verdhan, 2020). We used three steps to identify the best-performing classifier for the feature data (Figure 2): data processing, model optimization, and selection. We also compared the ability to classify unseen (test set) data using non-parametric, tree- and ensemble-based ML models. We found that the

**Table 1**  
Main Sample Characteristics, and Proportion of Main Particle Types in VolcAshDB

Samples	Eruption date	Magma composition	Volcano type	Eruptive activity type	Number of particles per component and associated error <sup>a</sup>					
					Altered material	Free-crystal	Juvenile	Lithic	Total	
Cumbre Vieja										
CV-DB1	19/10/21	Mafic	Cinder cone	Lava fountaining	3(±3)	1(±2)	719(±30)	352(±30)	1,075	
Kelud										
KE-DB2	14/2/14	Intermediate	Stratovolcano	Subplinian	50(±13)	4(±4)	268(±13)	3(±3)	325	
KE-DB3	14/2/14	Intermediate	Stratovolcano	Subplinian	162(±18)	59(±14)	54(±13)	65(±14)	340	
Merapi										
ME-DB1	22/7/13	Intermediate	Stratovolcano	Dome explosion	232(±16)	13(±7)	0	78(±15)	323	
ME-DB2	22/11/13	Intermediate	Stratovolcano	Dome explosion	595(±23)	76(±16)	4(±4)	100(±18)	775	
Sourfière de Guadeloupe										
SG-DB1	8/7/76	Intermediate	Stratovolcano	Phreatic	222(±17)	54(±13)	0	66(±14)	342	
SG-DB2	1/3/77	Intermediate	Stratovolcano	Phreatic	134(±5)	8(±5)	0	0	142	
Nevados de Chillán										
NC-DB15	3/4/18	Intermediate	Dome complex	Dome explosion	224(±26)	77(±17)	92(±18)	749(±31)	1,142	
NC-DB2	29/12/16	Intermediate	Dome complex	Dome explosion	99(±16)	12(±7)	14(±7)	171(±17)	296	
Ontake										
ON-DB1	27/9/14	Intermediate	Stratovolcano	Phreatic	777(±0)	0	0	0	777	
Pinatubo										
PI-DB1	2/4/91	Silicic	Caldera	Phreatic	386(±19)	104(±18)	0	16(±8)	506	
Mount St Helens										
MS-DB1	18/5/80	Silicic	Stratovolcano	Plinian	4(±4)	0	255(±58)	2(±3)	261	
Total					2,888(±78)	408(±38)	1,406(±65)	1,602(±68)	6,304	

<sup>a</sup>The associated error is calculated using the equation of margin of error in Benet et al. (2024) at a confidence interval of 95% and expressed in absolute values.



**Figure 2.** Illustration of the steps involved from the dataset to the outcomes, including those to obtain the best optimized model, XGBoost. (1) Data processing of the full dataset (features and particle types), including the oversampling of the training set. (2) hyperparameter optimization and cross-validation to obtain the models with the highest cross-validation scores. (3) evaluation of the models with the test set (unseen by the model) and selection of XGBoost with the highest classification scores. The XGBoost classifier was applied for the prediction of particle types and feature importance. See more details in the main text and subsequent figures.

XGBoost model had the best scores as found in other fields, such as a scalable tree boosting for data mining (Chen & Guestrin, 2016) and intrusion detection systems (Dhaliwal et al., 2018). The XGBoost model was used to gain insights into the most important features by calculating the Shapley values (see Section 2.2.4 for their definition) and with feature permutation (Molnar, 2021). The steps mentioned above were automated with a Python program that was run using the Gekko cluster at the Nanyang Technological University (NTU) High Performance Computing Center.

### 2.2.1. Data Processing

The dataset consists of 33 features measured from each particle (Table S2 in Supporting Information S1) and the particle types (i.e., our target variable; Figure 2). The dataset is made of about 6,300 particles and was divided into a training set (80% of the total particles) to optimize and fit the models, and a test set (20%), not used during the model's learning process. The original features' distributions are heterogenous and were standardized using the Scikit-learn's function *StandardScaler*, and it is commonly done to ease the convergence of ML models (Géron, 2017). The standard scaler redistributes the values of each feature with the mean at 0, and the first standard deviation at 1 and -1. The features from the test set were also standardized according to the scaler that was fitted into the training set to avoid data leakage (Géron, 2017). Any outliers, defined as values higher and smaller than two standard deviations (Verdhan, 2020), were kept after visually confirming that the source images had no errors. Highly correlated variables were kept for estimating their importance for classification in the step of feature permutation (more details are reported in "Explaining the model's predictions" in Section 2.3.4). Highly correlated variables may cause multi-collinearity issues in regression models, but these have not been reported in tree-based models (Kotsiantis, 2013).

The VolcAshDB dataset contains more altered material than juvenile and lithic particle types, and free crystals are relatively scarce (Table 1). Such an uneven distribution of particle types may cause an imbalanced dataset problem. We addressed this issue by oversampling the less abundant particle types using the SMOTE package, which uses a K-Nearest Neighbor algorithm (KNN) to generate synthetic data (Brownlee, 2020). This step is strongly recommended to allow the model to learn to classify the less abundant class (Brownlee, 2020).

### 2.2.2. Hyperparameter Optimization

Hyperparameters are user-defined settings that control the models' learning process. In our study, we explored various models including Decision Trees (DTC), K-Nearest Neighbor (KNN), Random Forest (RF), Gradient Boost Classifier (GBC), and the Extreme Gradient Boosting (XGBoost). To find the best hyperparameter values

more efficiently, we used Bayesian optimization through the function *BayesSearchCV* from the Scikit-optimize package. This function searches for the optimal hyperparameters depending on the previous iterations, making the computation faster and less intensive than iterating through the entire search space (Owen, 2022). To evaluate the effect of the different values of hyperparameters, we obtained scores from 10-fold cross-validation on the training set. This method divides the data into training and testing folds iteratively and is recommended to avoid overfitting (Verdhan, 2020). Using the optimal hyperparameters on the tested models (Table S4 in Supporting Information S1), XGBoost yielded the highest cross-validation score with 0.9 *F1-score* (a metric defined and calculated below in Section 2.2.3), closely followed by KNN and GBC with 0.88 *F1-score* (obtained scores of each model are shown in Figure S1 in Supporting Information S1).

### 2.2.3. Model Evaluation and Selection

The cross-validation scores indicate how well a model fits the training set. To evaluate the models' ability to generalize, we also computed the predictions on the test set. Each prediction contains a confidence score per class, which represents the likelihood of the prediction belonging to the class, and the score is given as a percentage (Mandal et al., 2021). The class, that is, the particle type in this study, with the highest confidence score is considered the predicted type by the model. Comparison between the predicted and the visually assessed types from VolcAshDB allows to categorize each prediction in one of the four following groups: True Positive (TP), where the prediction correctly identifies the class; True Negative (TN), where the prediction correctly identifies the absence of a class; False Positive (FP), where the prediction wrongly identifies the presence of a class, and False Negatives (FN), where the prediction wrongly identifies the absence of a class. The classification matrix (Figure S2 in Supporting Information S1) is typically used in ML to show the proportions of TP, TN, FP and FN for each class. Based on these proportions, we can calculate four well-known metrics to evaluate the models' performance (e.g., Verdhan, 2020):

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}} \quad (1)$$

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}} \quad (2)$$

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}} \quad (3)$$

$$F1 - \text{score} = \frac{2 * \text{TP}}{2 * \text{TP} + \text{FP} + \text{FN}} \quad (4)$$

We report the classification performance using the *F1-score*, as it combines the precision, dependent on the FP, and recall, dependent on the FN, into a single metric (Verdhan, 2020) and is recommended for imbalanced datasets when FN and FP are equally important (Brownlee, 2020). We use the unweighted average of the *F1-scores* (the so-called *macro* from macro-averaging) of the four particle types to evaluate the overall model performance, as opposed to the weighted averaging, where the average is multiplied to a coefficient based on the number of particles per class (Verdhan, 2020). We found that XGBoost has the best classification performance with 0.77 *macro F1-score* amongst the optimized models and therefore is our selected model (classification score for each model are reported in Table S5 in Supporting Information S1 and shown in Figure S3 in Supporting Information S1).

### 2.2.4. Explaining the Model's Predictions

Explainable Artificial Intelligence (xAI) is a set of methods that provide explanations on the variables that drive the model's predictions (Gianfagna & Di Cecco, 2021; Mishra, 2022; Molnar, 2021). We gained insights into the contribution of the 33 features to the model's predictions by applying two methods: "permutation feature importance" and the SHapley Additive exPlanations (SHAP; Lundberg & Lee, 2017). Permutation feature importance involves shuffling the values of each feature to measure the resulting increase in prediction error. The features are then ranked based on their contribution to prediction error. "Important" features are those causing an increase in error when shuffled, whereas "unimportant" features exhibit no change or a decrease in error (Molnar, 2021).

To implement this method, we used Scikit-learn's *permutation* function on the test set. We estimated the feature importance of each class by permuting the features between each class and the rest (e.g., One-vs-Rest strategy).

The SHAP method, implemented as a Python library (Lundberg & Lee, 2017), has been key for explaining individual model's predictions in both regression (e.g., Biass et al., 2022; Kondylatos et al., 2022) and classification problems (e.g., Panati et al., 2022; Tang et al., 2021). The method is based on the Shapley values, a concept introduced by Shapley (1953). These values measure the contribution of the feature values in predicting a specific outcome (e.g., a particle type in our case) relative to the average prediction across all instances (Molnar, 2021). To calculate the Shapley values, we used the TreeSHAP estimation method with raw output. Because Shapley values are additive, the TreeSHAP method adds and averages the contribution of each node in the ensemble trees to obtain the Shapley value of each feature value per instance (Lundberg et al., 2018). The highest positive Shapley values per instance are those which contribute the most to the prediction of a given particle type. By averaging Shapley values by particle type or across all four types (altered material, free-crystal, juvenile, and lithic), we gained insights into global feature importance (Lundberg et al., 2018), which can be used for comparison with the permutation feature importance.

### 2.2.5. Classification Strategies

We applied three classification strategies to evaluate which model performs best: (a) the multilabel, where the four classes are used to train the model at once, with a prediction probability given for each class, and that with the highest value is the predicted class, (b) the One-vs-One (OVO), where each possible pair of classes trains a binary classifier (i.e., a total of six classifiers, as there are six possible pairs for four classes), and their outputs are aggregated to yield the predicted class (Herrera et al., 2016), and (c) the One-vs-Rest (OVR), where each class and its complementary (e.g., lithic vs. non-lithic) are used to train a binary classifier (i.e., a total of four), and their outputs are aggregated to yield the predicted class (Herrera et al., 2016). For the OVO and OVR strategies, the outputs from the binary classifiers were aggregated with the same weight to obtain the predicted class. There are more sophisticated aggregation methods, such as the calibrated label ranking method (Fürnkranz et al., 2008), which adjust the weights of each binary classifier aiming to mitigate class dependencies and make the global classification more robust (Herrera et al., 2016). However, we do not know of any implementation of these methods in Python for the XGBoost model, and developing them from scratch is out of the scope of this study.

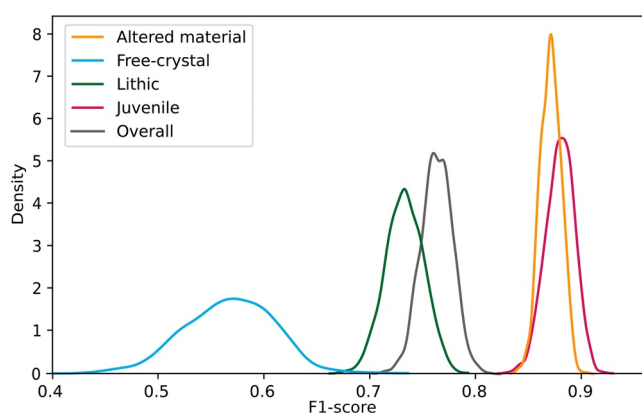
### 2.2.6. Effect of the Training and Test Data Split on the XGBoost Scores

As noted above, we first split the dataset into a training (80% of all particle features in VolcAshDB) and a test set (20%) and used the latter to evaluate the XGBoost's performance. However, as the splitting process is random, it may affect the precision and accuracy of the measured *F1-scores*. To estimate this error, we re-trained and evaluated the XGBoost at 1,000 different values of random state, that is, the hyperparameter that controls randomness. We obtained an average accuracy (*macro F1-score* of 0.76; Table S6 in Supporting Information S1) that is similar to the accuracy from the test set (*macro F1-score* of 0.77). The free-crystal type shows the widest variability (standard deviation of 0.04) and is the most inaccurate (*F1-score* of 0.57; Figure 3) among the particle types. This is likely because it is the least abundant type, and its classification is challenging given the different types of minerals and lack of a discriminant feature, as explained (Section 3.1). Accuracies of the three other types are higher (*F1-score* of 0.73–0.88) and with better precision (standard deviation is <0.02; Table S6 in Supporting Information S1).

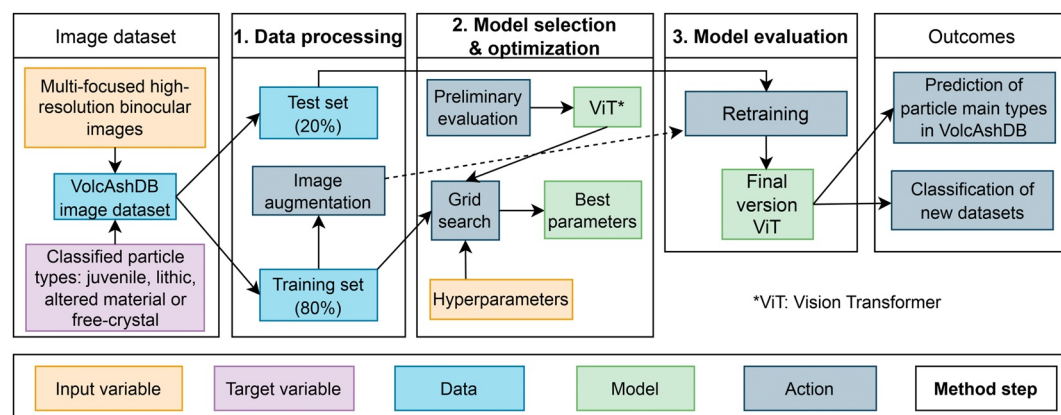
By averaging the *F1-scores* of each particle type, we obtain the *macro F1-score* distribution (Figure 3) and its variability (standard deviation; Table S6 in Supporting Information S1). To quantify the associated error ( $\alpha$ ), we use the second standard deviation (Hughes & Hase, 2010):

$$\alpha = \left[ \sqrt{\frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{x})^2} \right] * 2 \quad (5)$$

where  $N$  is the number of experiments,  $x$  is each measured value (i.e., *macro F1-score*) and  $\bar{x}$  is the mean. With the values noted above, we obtain an error ( $\alpha$ ) of 0.03 for *macro F1-score* distribution, and since we used the second



**Figure 3.** Density plots of the *F1-scores* obtained from 1,000 runs of the XGBoost at different random states across particle types and aggregated as *macro F1-score* (Overall).



**Figure 4.** Illustration of the steps involved from the dataset to the outcomes, including those to fine-tune the Vision Transformer (ViT). (1) Data processing of the full dataset (images and particle types). (2) preliminary evaluation of the models using the base hyperparameters, selection of ViT and hyperparameter optimization through grid search. (3) Final evaluation using the test set after fine-tuning with the augmented dataset. The ViT classifier can be then applied for the prediction of particle types. See more details in the main text and subsequent figures.

standard deviation, it is for a 95% confidence level. Therefore, the *FI-score* values can be reported as:  $0.76 \pm 0.03$  *macro FI-score*, which is a small relative error of <5%.

### 2.3. Development of a Particle Classifier Using VolcAshDB Image Data Set

We used three steps to develop an optimized classifier for the image dataset (Figure 4): data processing, model selection and optimization, and model evaluation. We compared the performance between three state-of-the-art models that have top accuracies in the reference dataset ImageNet (J. Deng et al., 2009): ResNet (He et al., 2016), which is the prevalent model of the so-called convolutional neural networks (CNN), Vision transformer (ViT; Dosovitskiy et al., 2020), which superseded ResNet in image classification and ConvNeXT (Z. Liu et al., 2022), which is an optimized convolutional neural network that has surpassed the performance of vision transformers. The models are available in the *Hugging Face* library (<https://huggingface.co/>), which also provides application programming interfaces (API) for their deployment. The model that yielded the highest classification score was the ViT. We augmented the training dataset with an array of variations from the original images (see below), and the ViT reached a *macro FI-score* of 0.93, outperforming the XGBoost classifier (Table 2). The images of the ash particles in VolcAshDB were obtained from processed multi-focused binocular images, but this is not the standard practice, and thus we also tested the ViT's ability to classify standard single-focus binocular images used in most studies. The steps mentioned above were automated with a Python program that was run using the Gekko cluster at the Nanyang Technological University (NTU) High Performance Computing Center.

#### 2.3.1. Image Augmentation and Processing

The binocular images of ash particles in VolcAshDB are multi-focused and contain four color-related channels: red, green, blue, and alpha. The alpha channel is a binary mask that takes a value of 1 or 0 to separate between the particle pixels and those of the background (more details in the segmentation step in Benet et al. (2024)). We split the dataset into a train (80% of the total images in VolcAshDB) and test set. Then, we augmented the number of images in the training set based on six standard methods (Ayyadevara & Reddy, 2020): rotation (at 45°), translation (at 25 pixels), up-down and left-right flipping, and adding random noise and Gaussian blur at sigma values of 0.155 and 0.55. Increasing the number of images allowed us to balance the distribution across particle types (~2,900/class), and is generally recommended to increase model's robustness (Brownlee, 2020). Images were

**Table 2**  
*FI-Score Values for the Whole Database (Unweighted Average or Macro) and Particle Types Obtained From Various Models, Including XGBoost Multilabel, One-vs-One (OVO), One-vs-Rest (OVR), and the Multilabel Image-Based Model ViT*

	Overall	Altered material	Free-crystal	Juvenile	Lithic
Multilabel XGBoost	0.77	0.88	0.57	0.90	0.74
OVO XGBoost	0.75	0.89	0.56	0.85	0.71
OVR XGBoost	0.76	0.90	0.55	0.88	0.73
Multilabel ViT	0.93	0.95	0.91	0.95	0.89

stored in four subdirectories, one for each class, of a root directory which was inputted to the *Hugging Face's API* for fine-tuning.

### 2.3.2. Fine-Tuning, Preliminary Evaluation, and Model Selection

We fine-tuned the classifiers and performed a preliminary round of evaluations to choose the best-performing model. Fine-tuning consists of making small adjustments to an already trained classifier, as opposed to training, where the data drives the model's learning process without any prior exposure. We selected the model before hyperparameter optimization because each run is time consuming (lasting about 14–18 hr) and because the authors of each model already provide the base hyperparameters (Table S7 in Supporting Information S1). The fine-tuned model that yielded the highest accuracy is ViT (0.88), followed by ConvNext and ResNet, both with an accuracy of 0.86.

### 2.3.3. ViT Hyperparameter Optimization

We obtained the optimal hyperparameters following the grid search technique for two ranges of batch size and learning rate. In grid search, each hyperparameter is modified one step at a time, while the other hyperparameters remain fixed throughout all the possible combinations (Owen, 2022). We found that multiple combinations yielded similar accuracies exceeding 0.87 (accuracies obtained from grid search are reported in Table S8 in Supporting Information S1) and opted for the batch size and learning rate at 16 and  $3 \times 10^{-5}$ , as these are more commonly used according to established practices in deep learning (Kandel & Castelli, 2020). Using these values, we tested three different optimizers, AdamW (Loshchilov & Hutter, 2019), Stochastic Gradient Descent (Sutskever et al., 2013) and Adagrad (Duchi et al., 2011), with the former performing the best (Table S9 in Supporting Information S1). We also tested an increasing number of epochs (i.e., 5, 10, 15, 20), which did not improve performance above 10, probably because the model had already converged.

### 2.3.4. Model Evaluation

We fine-tuned again the ViT with the augmented training set and the optimal set of hyperparameters and obtained a significant improvement, with a *macro F1-score* of 0.93. We obtained the same metrics of precision, recall, accuracy and F1-score, confusion matrix, and confidence scores as defined and calculated above (Section 2.2.3 Model evaluation and selection). In contrast with the XGBoost, the explainability of the model is very limited, as further discussed below (see Section 4.1).

## 3. Results

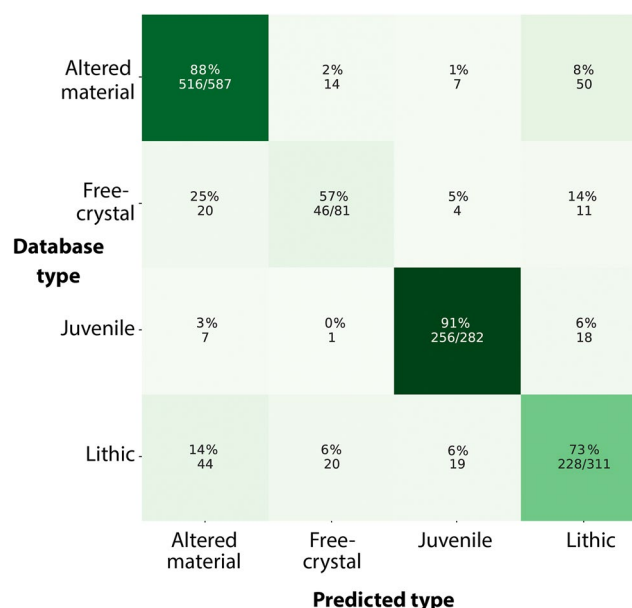
We used the VolcAshDB ash particle features and images to train the XGBoost and ViT models and to evaluate their ability to classify them into altered material, free-crystal, juvenile, or lithic types (Table 2). We found that overall ViT classifies very accurately, with a *macro F1-score* of 0.93, whereas the XGBoost is less performant with a *macro F1-score* of 0.77 (Table 2) but allows for explaining the model's predictions using interpretable AI methods. We describe below the model performance through the two datasets by particle type and some particle subgroups, such as those divided by the volcano, or one class versus another.

### 3.1. XGBoost Quantitative Evaluation

Overall, the XGBoost shows rather accurate *F1-scores* across classification strategies: 0.77 for multilabel, 0.75 for OVO, and 0.76 for OVR (Table 2). The computation of the confusion matrix (Figure 5) shows that the model classifies the best juvenile type (*F1-score* of 0.9), closely followed by altered material (*F1-score* of 0.88), and less accurately the lithic type (*F1-score* of 0.74), and significantly less the free-crystal type (*F1-score* of 0.57).

Binary classifications using OVO and OVR between altered material, lithic and juvenile have accuracies  $>0.80$  (*macro F1-scores* of 0.82–0.97), whereas the free-crystal type is systematically lower (Table S10 in Supporting Information S1). A closer inspection by volcano and eruptive activity reveals a wide range in XGBoost's performances (Table 3). Predictions of juvenile particles are very accurate (*F1-score* of 0.97) at Kelud volcano but inaccurate (*F1-score* of 0.32) at Nevados de Chillán. Classification of lithics is rather accurate for samples of dome explosions (*F1-score* of 0.8) but inaccurate (*F1-score* of 0.28) for those of phreatic events. Such fluctuations indicate limited robustness by the classifier and care should be taken for its application to other datasets on a case-by-case basis.

The likelihood that a particle belongs to a given type according to the model is reflected in the distribution of the confidence scores and varies across particle types. Within the True Positives (TP) group, almost 90% of the



**Figure 5.** Confusion matrix of the predictions by the XGBoost multilabel classifier. The percentages show the True Positive rate if positioned in the diagonal matrix (darker green), and otherwise the False-Negative rate (lighter), all percentages with the corresponding number of particles per predicted type. The best classification is for juvenile followed in descending order by altered material, lithic and free-crystal types.

juvenile TP have confidence scores  $>0.9$ , whereas  $\sim 40\%$  of the free-crystal TP have confidence scores between 0.4 and 0.9 (Figure 6a). This means that the XGBoost is almost certain when predicting juvenile particles, but more unstable for free crystals. The confidence scores over the False Negatives (FN) show that the XGBoost identifies a relatively high number of lithic particles and free crystals as altered material, with confidence scores  $>0.9$  (Figures 6b and 6c), hinting at some classification challenges that are revealed below using the Shapley values (see “Local feature importance” in Section 3.2.2).

### 3.2. What Features Drive XGBoost Particle Type Predictions?

#### 3.2.1. Global Feature Importance

We identified the features driving the XGBoost's predictions with two approaches: applying the permutation feature importance and computing the mean of the Shapley values (see Section 2.3.4). Although the calculation of the two methods is quite different, they yielded overall a similar feature importance ranking. We identified the following three as the most important features (Table 4): (a) the mean of the hue channel

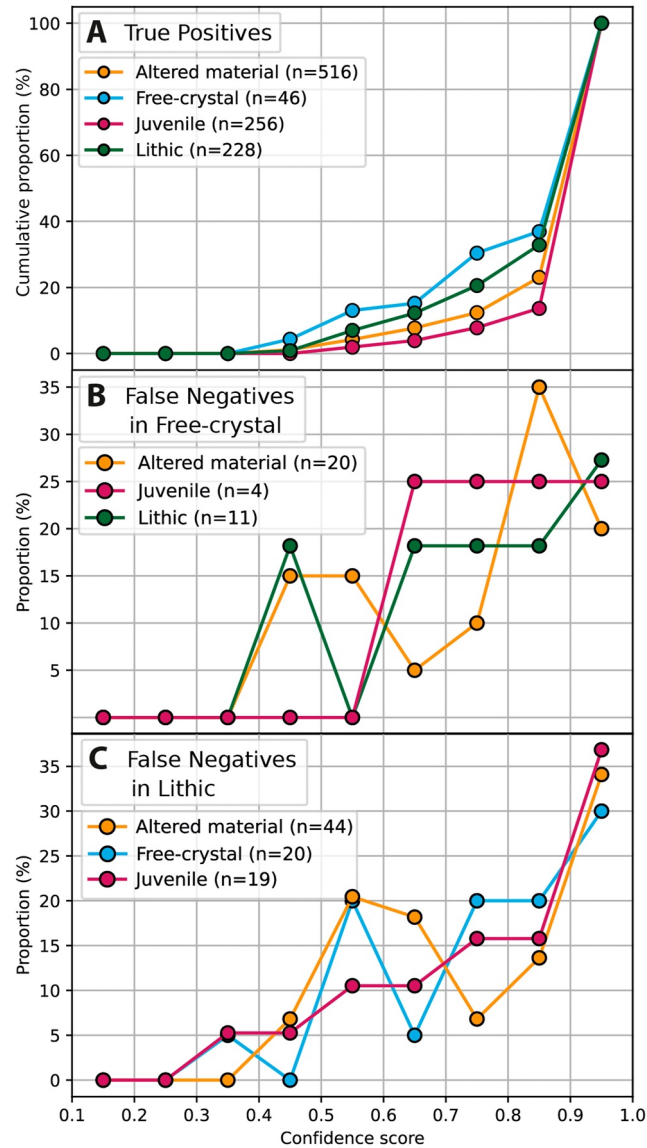
**Table 3**

*F1-Scores Obtained From the Multilabel XGBoost Classifier of Each Particle Type and Their Unweighted Average (i.e., Macro) for All Particles in the Test Set (Overall), and Across Volcanoes and Eruptive Activity Types*

	Volcano					Eruptive activity				
	Overall	Soufrière de Guadeloupe	Merapi	Nevados de Chillán	Cumbre Vieja	Kelud	Phreatic	Dome explosion	Lava fountaining	Sub-plinian/Plinian
<i>F1-score (macro)</i>	0.77	0.76	0.73	0.6	0.87	0.73	0.62	0.65	0.87	0.76
<i>A<sup>a</sup></i>	0.88	0.92	0.91	0.7	–	0.81	0.95	0.82	–	0.84
<i>F<sup>b</sup></i>	0.57	0.7	0.67	0.59	–	0.6	0.64	0.51	–	0.7
<i>J<sup>c</sup></i>	0.9	–	–	0.32	0.92	0.97	–	0.46	0.92	0.99
<i>L<sup>d</sup></i>	0.74	0.67	0.6	0.77	0.83	0.54	0.28	0.8	0.83	0.42

*Note.* These measurements have an estimated precision of  $\pm 0.03$ .

<sup>a</sup>A: Altered material. <sup>b</sup>F: Free-crystal. <sup>c</sup>J: Juvenile. <sup>d</sup>L: Lithic.



**Figure 6.** Line plots of the confidence score versus (a) the cumulative proportion of True Positives (TP), (b) False Negatives (FN) in free-crystals, and (c) lithic types. The distribution of the data has been plotted into 9 bins of size 0.1. We do not use cumulative proportion in (b, c) given the limited number of FNs. The meaning of the plot in (a) can be understood by the following two examples: if we take the value of juvenile TP at a confidence score between 0.8 and 0.9, there is a low cumulative proportion of ~10%, whereas in the next bin, 0.9–1.0 of the confidence score, we have the vast majority (~90%) of the juvenile TP. If we take the value of free-crystal TP at a confidence score between 0.8 and 0.9, there is a significant cumulative proportion of almost 40%. This shows that XGBoost is more reliant on predicting juvenile particles than free crystals.

(*hue\_mean*), which is a feature from the Hue-Saturation-Value color space that measures the averaged chromaticity; (b) the *correlation*, a textural feature that measures the degree of similarity between pixel relationships (Hall-Beyer, 2017); and (c) the mode of the blue channel (*blue\_mode*), which measures the most frequent pixel intensity of the blue channel of the particle image.

### 3.2.2. Local Feature Importance Across Particle Types

We identified the most important features used by the XGBoost to predict each particle type based on the Shapley values, which consider the interaction between the four particle types, unlike permutation, which is based on the One-vs-Rest approach. Shapley values calculate the contribution of each feature to the actual prediction with respect to the expected prediction (Gianfagna & Di Cecco, 2021; Lundberg et al., 2018; Molnar, 2021). Thus, we can use the Shapley values of an individual particle prediction to identify which features were more

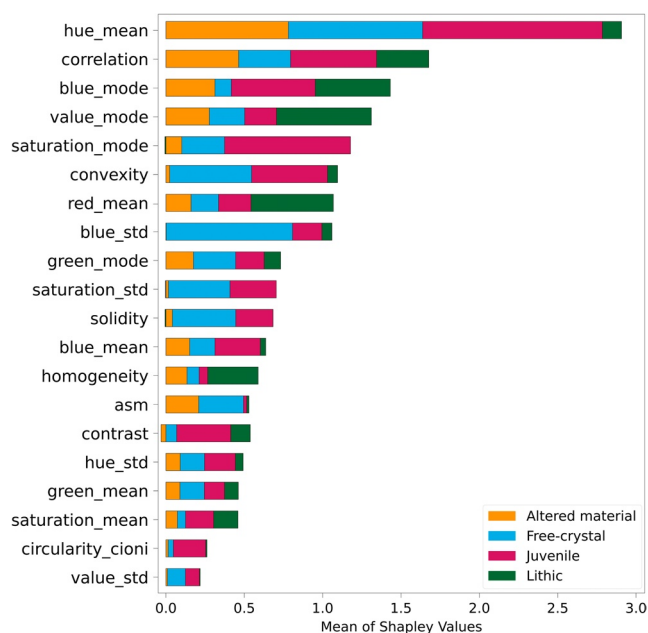
**Table 4**  
Feature Importance Identification Based on Mean of Shapley Values and Feature Permutation

Feature importance method <sup>a</sup>	Mean of shapley values					Feature permutation				
	Per particle type (multilabel)					Per particle type (OVR)				
	A	F	J	L	Total	A	F	J	L	Total
<i>hue_mean</i>	<b>0.78</b>	<b>0.86</b>	<b>1.15</b>	<b>0.12</b>	<b>2.91</b>	<b>0.91</b>	<b>0.41</b>	<b>0.91</b>	<b>0.15</b>	<b>1.22</b>
<i>correlation</i>	<b>0.46</b>	<b>0.33</b>	<b>0.55</b>	<b>0.33</b>	<b>1.68</b>	<b>0.34</b>	0.02	<b>0.04</b>	<b>0.19</b>	<b>0.29</b>
<i>blue_mode</i>	<b>0.31</b>	0.10	<b>0.54</b>	<b>0.48</b>	<b>1.43</b>	0.06	<b>0.04</b>	<b>0.01</b>	0.00	<b>0.10</b>
<i>value_mode</i>	<b>0.28</b>	<b>0.23</b>	0.20	<b>0.60</b>	<b>1.31</b>	0.05	<b>0.05</b>	0.00	<b>0.24</b>	0.00
<i>saturation_mode</i>	<b>0.10</b>	<b>0.27</b>	<b>0.80</b>	−0.01	<b>1.17</b>	0.02	<b>0.06</b>	<b>0.10</b>	<b>0.10</b>	<b>0.13</b>
<i>convexity</i>	0.02	<b>0.52</b>	<b>0.48</b>	0.06	<b>1.10</b>	0.01	<b>0.06</b>	<b>0.03</b>	0.00	0.03
<i>red_mean</i>	<b>0.16</b>	0.18	<b>0.21</b>	<b>0.53</b>	<b>1.07</b>	0.03	0.03	<b>0.01</b>	0.01	0.04
<i>blue_std</i>	−0.06	<b>0.81</b>	0.19	0.06	<b>1.00</b>	<b>0.34</b>	<b>0.24</b>	<b>0.04</b>	<b>0.04</b>	<b>0.28</b>
<i>green_mode</i>	<b>0.18</b>	<b>0.27</b>	0.18	<b>0.11</b>	<b>0.73</b>	0.03	0.02	<b>0.03</b>	0.01	0.02
<i>saturation_std</i>	0.02	<b>0.39</b>	<b>0.30</b>	0.00	<b>0.70</b>	<b>0.07</b>	0.00	<b>0.08</b>	0.00	<b>0.11</b>
<i>solidity</i>	0.04	<b>0.40</b>	<b>0.24</b>	−0.01	0.68	<b>0.08</b>	0.01	<b>0.02</b>	<b>0.07</b>	−0.04
<i>blue_mean</i>	<b>0.15</b>	0.16	<b>0.29</b>	0.03	0.64	0.06	<b>0.05</b>	<b>0.01</b>	0.01	0.05
<i>homogeneity</i>	<b>0.13</b>	0.08	0.06	<b>0.32</b>	0.59	<b>0.16</b>	0.03	0.00	<b>0.12</b>	<b>0.06</b>
<i>asm</i>	<b>0.21</b>	<b>0.29</b>	0.02	0.01	0.53	<b>0.18</b>	0.03	0.00	0.00	<b>0.14</b>
<i>contrast</i>	−0.03	0.07	<b>0.35</b>	<b>0.12</b>	0.51	<b>0.11</b>	0.03	0.00	0.02	0.03
<i>hue_std</i>	0.09	0.16	0.20	0.05	0.49	<b>0.14</b>	<b>0.13</b>	0.00	<b>0.11</b>	<b>0.14</b>
<i>green_mean</i>	0.09	0.16	0.13	<b>0.09</b>	0.46	<b>0.16</b>	0.02	0.00	<b>0.13</b>	<b>0.13</b>
<i>saturation_mean</i>	0.07	0.05	0.18	<b>0.15</b>	0.46	0.01	<b>0.05</b>	<b>0.01</b>	0.00	0.04
<i>circ_cioni</i>	0.01	0.03	0.21	0.01	0.26	0.01	0.00	−0.01	0.02	−0.02
<i>energy</i>	0.05	0.02	0.00	0.06	0.14	0.03	0.00	0.00	<b>0.09</b>	0.01
<i>red_std</i>	−0.01	0.00	0.09	0.03	0.11	0.03	<b>0.13</b>	0.00	0.00	0.03
Total	3.12	5.51	6.51	3.13		2.86	1.43	1.29	1.33	

<sup>a</sup>These two methods calculate the feature importance values differently and cannot be directly compared. The relative ranking of the feature importance is similar (top 10 ranked features in bold) to the same top two ranked features (*hue\_mean* and *correlation*). We used the Shapley mean value for feature importance per particle type (shown as a plot in Figure 7), the top three of which are underlined. For the meaning of the abbreviations of each feature please see Table S2 in Supporting Information S1. The permutation feature values have been multiplied by 10 for better readability, as the importance lies on the relative values across features.

important, or average them across particle types to identify the global discriminant features per type (Figure 7) as follows:

1. *Altered material* has the second highest classification performance with a *FI-score* of 0.88 and is predicted with color (*hue\_mean*) and texture (*correlation*) (Figure 8a). A group of True Positives (*TP*) with *hue\_mean* values between −3 and −2 (rescaled as described in Section 2.3.1) is revealed by the Shapley dependence plot (Figure 8b), which relates feature values (*hue\_mean*) and their associated Shapley values for each particle (Lundberg et al., 2018). Such *TP* have almost 100% of confidence scores and consist of white (Figure 8c), red (Figure 8d), rounded, hydrothermally altered material.
2. *Free crystals* are the least accurately classified with *FI-score* of 0.54 and are mainly discriminated by color (*hue\_mean* and *blue\_std*), and shape (*convexity*) (Figure 9a). Unlike the other types, the most discriminant feature does not cluster particles, as shown by the *blue\_std* values as a function of the Shapley values does not yield any cluster of *TP* (Figure 9b) and those with Shapley values > 1.5 overlap with altered material (Figure 9c). Thus, the XGBoost has limited predictability of free crystals, which is consistent with a low *FI-score* yielded from Free-crystal type versus Rest binary classification (Table S10 in Supporting Information S1). Possible causes for this, besides the lack of a discriminant feature, include the presence of glass films on the crystal's surface, the wide range of aspects of different minerals (mostly plagioclase and pyroxene



**Figure 7.** Aggregated mean of the Shapley values by particle type. Note that some features are important for the discrimination of multiple particle types (e.g., *hue\_mean*) and other features are more discriminant of a specific type (e.g., *value\_mode* for lithic type). The meaning of the abbreviations can be found in Table S2 in Supporting Information S1.

but also amphibole and sulfur-group minerals), and the significant rate of composite particles (e.g., crystals attached to glass) that are not reflected in the label (Figure 9d).

3. The *juvenile* particles have the highest classification success with an *F1-score* of 0.90 and are predicted through color (*hue\_mean*, *saturation\_mode*), texture (*correlation*), and shape (*convexity*) (Figure 10a). The *saturation\_mode* feature, which relates to the intensity of color, is discriminant (Shapley values > 1) with values of 0–2 (Figure 10b). Low values of *convexity* are also relevant for discrimination, as could be expected by the presence of vesicles on the particles' surfaces (Figure 10c). Moreover, the XGBoost predicts instances with lower *hue\_mean* and *saturation\_mode* as lithic (i.e., False Negative, FN), which correspond to darker, mid to high crystallinity juvenile particles from dome explosions (Figure 10d).
4. The *lithic* particles are moderately well classified with a *F1-score* of 0.74, and are mainly discriminated through color (*value\_mode*, *red\_mean* and *blue\_mode*) and texture (*correlation* and *homogeneity*) (Figure 11a). Low values of *value\_mode*, ranging between of  $-1.7$  to  $0$  (Figure 11b), are characteristic of a duller luster and discriminate lithic particles (Shapley values > 0.5). This shows that a dull luster may allow discriminating lithic particles from the rest, for example, glossy juvenile particles, as proposed by Miwa et al. (2009). These features together with relatively high values of *correlation*, reflect dark, dull lithic particles with uniform texture (Figure 11c). In contrast, instances with higher pixel intensity-based features (*green\_mean*) are a source of FN, as suggested by the negative Shapley values, and are classified as altered material (Figure 11d).

### 3.3. ViT Quantitative Evaluation

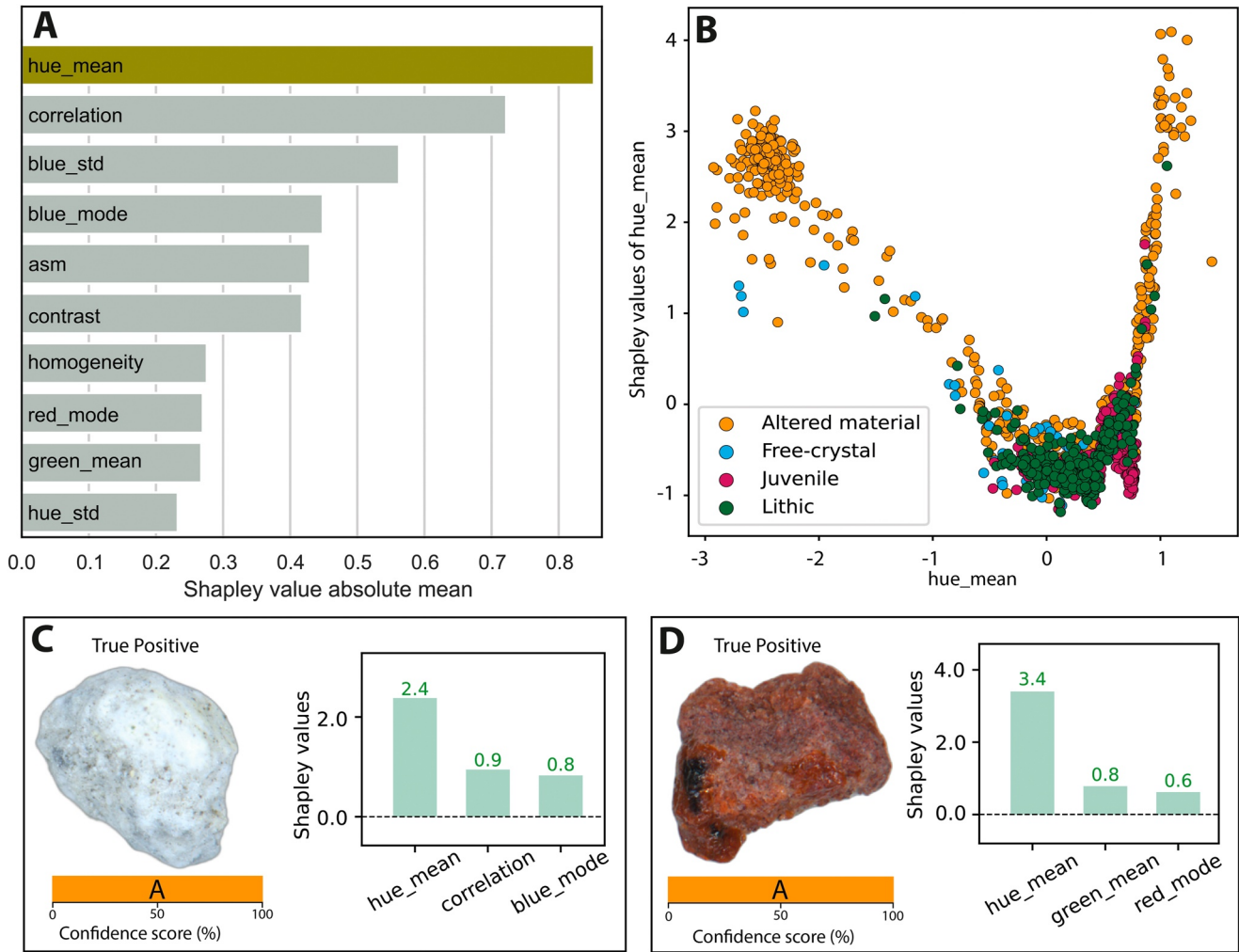
#### 3.3.1. General Evaluation

The ViT base model was fine-tuned using  $\sim 10,000$  images from the augmented training set and evaluated with the test set (see Section 2.3 for information on each step). We obtained an accurate classification for the whole test set (*macro F1-score* of 0.93), and also across particle types (Figure 12): altered material (*F1-score* of 0.95), juvenile (*F1-score* of 0.95), free-crystal (*F1-score* of 0.91) and lithic (*F1-score* of 0.89). More than 85% of True Positives (*TP*) are predicted at high confidence scores (>0.9; Figure 13a) which shows that ViT classifies confidently and accurately. The False Negatives (*FN*) mostly consist of lithic particles classified as altered material and juvenile, a few of which at high confidence scores (Figure 13b), and also of juvenile particles classified as lithic type (Figure 13c). Below, we identify specific groups of particles that make up the *FN* and discuss the possible causes.

#### 3.3.2. ViT's Evaluation Across Volcanoes, Eruptive Activity Types, and Individual Particles

A closer inspection of the results across eruptive activity types and volcanoes (Table S11 in Supporting Information S1) reveals a range of classification accuracies, from moderate (*F1-score* of 0.73) up to optimal classification performance with a *F1-score* of 1.0 (Figure 14):

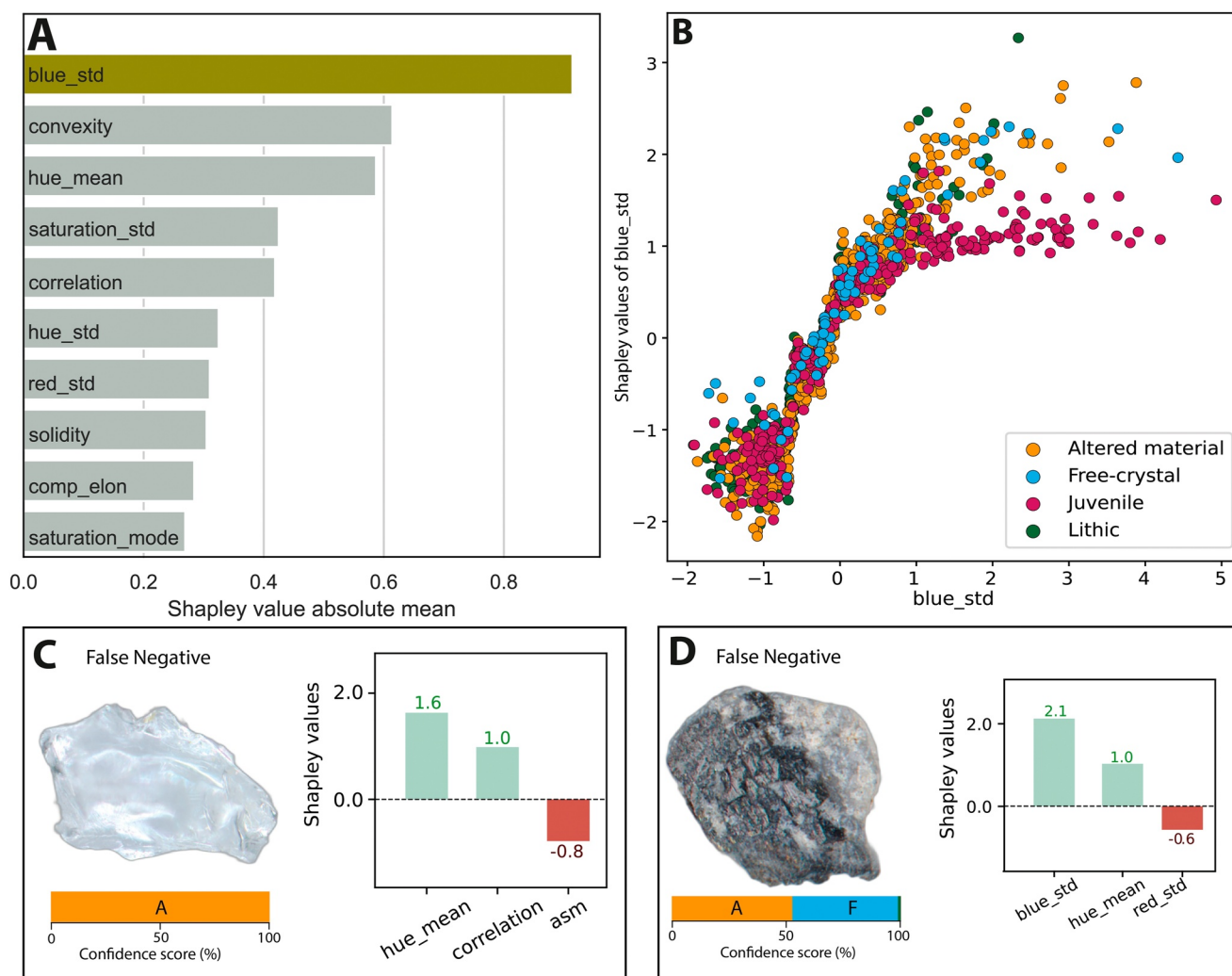
1. Ash particles from phreatic events are in general well classified (*macro F1-score* of 0.95), including the particle main types: altered material (*F1-score* of 0.99), free-crystal (*F1-score* of 0.94) and lithic (*F1-score* of 0.93). The ViT successfully classifies the most common groups of particles in these samples such as hydrothermal aggregates (Figure 15a) and weathered material (Figure 15b).
2. Particles from samples of dome explosions are classified with the lowest accuracy (*macro F1-score* of 0.85) among the eruptive activity types. The ViT accurately classifies free-crystal (*F1-score* of 0.86), altered material (*F1-score* of 0.90) and lithic (*F1-score* of 0.90) types, but is less accurate (*F1-score* of 0.73) for the juvenile type with most False Negatives (FN) classified as lithics. However, the confidence scores of some FN



**Figure 8.** Summary plots to explain the predictions of the altered material particle type. (a) Feature importance according to the mean of the Shapley values; the higher the value, the more important the feature in the correct prediction. In (b), the Shapley dependence plot shows the relation of the Shapley value and the feature value for each particle type, and is commonly used to identify clusters of a specific class (particle main type) along the feature domain (Lundberg et al., 2018). For example, at values of  $-3$  to  $-2$  of *hue\_mean*, there is a cluster of particles with high Shapley values and thus correctly classified as altered material. (c, d) are two examples of particles to show confidence score (A: Altered material), and the three features with the highest Shapley values. They are both True Positives and have been predicted at maximum confidence score with *hue\_mean* (the mean of the chromaticity) being the main discriminant feature.

show a transition between the juvenile and lithic types that has an explanatory value. This means that particles may have both juvenile and lithic traits, and thus a measure on the types' prevalence seems more realistic than using mutually exclusive types like in VolcAshDB. Particles with combined traits are common in samples from the Nevados de Chillán Volcanic Complex (Figure 15c), which originated from a relatively long-lived dome-forming eruption cycle. An additional challenge is that the ViT confidently classifies as lithics some particles that are labeled as juvenile, and since operator-based classification was not always straightforward (Benet et al., 2024), it is difficult to decide whether these are False Negatives or instead petrographic classification errors (Figure 15d), especially when ML-based image classifiers have surpassed human performances in other fields (He et al., 2015).

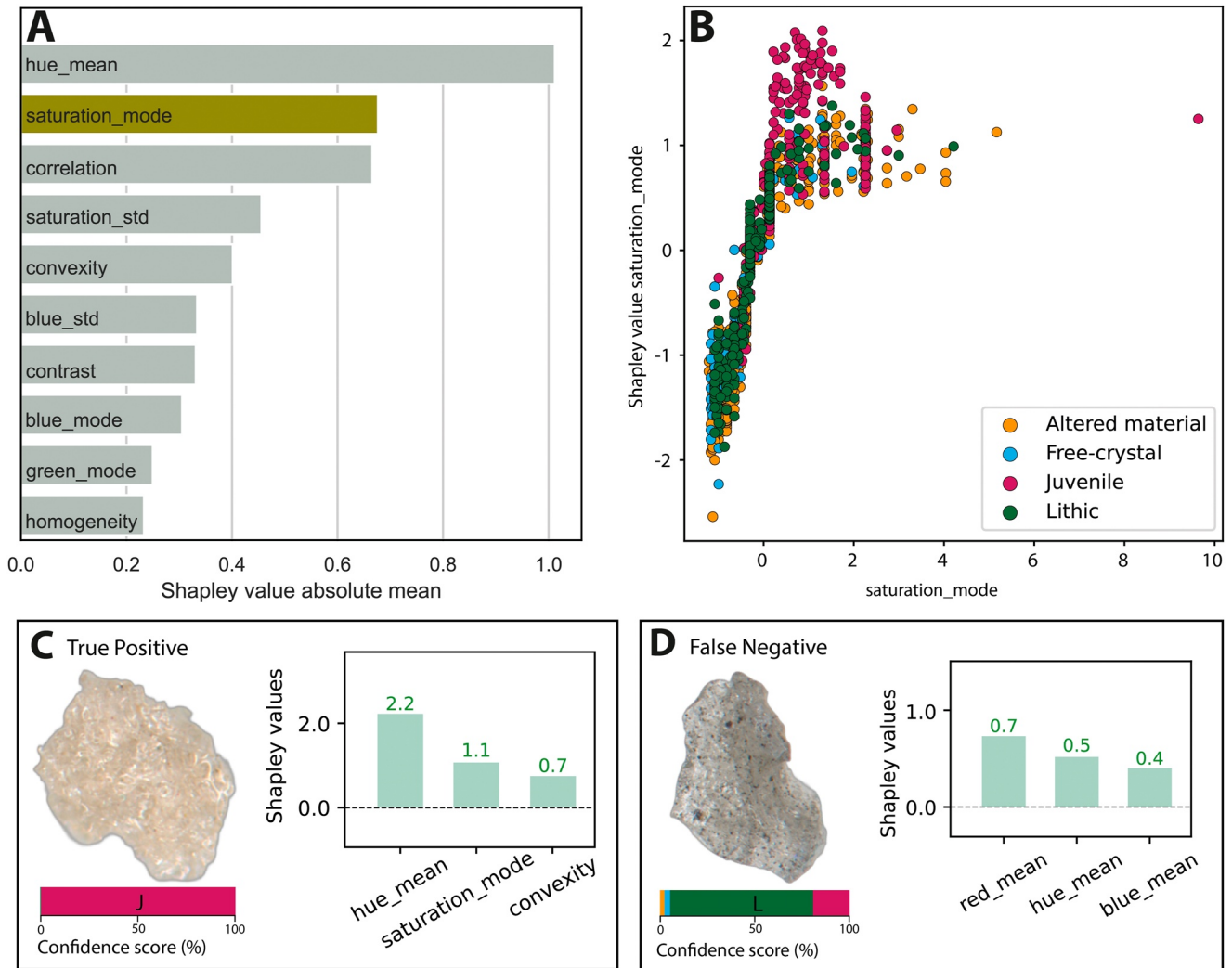
- Ash particles from lava fountaining are generally accurately classified (*macro F1-score* of 0.94) between juvenile (*F1-score* of 0.94) and lithic (*F1-score* of 0.88) types. Most of the lithic particles belong to recycled



**Figure 9.** Summary plots to explain the predictions of the models for the free-crystal type. (a) Feature importance based on the mean of the Shapley values. (b) Shapley dependence plot. Note that the feature values have been rescaled by a standard scaler. (c) and (d) show for each prediction: the particle image, the confidence score across particle types, and the associated Shapley values. (c) Shows a particle that is likely a fragment of plagioclase crystal but is misclassified as altered material, because the free-crystal type lacks discriminant features (see main text for more details). (d) An additional source of false negatives are particles consisting of more than one material, such as those made of glass attached to a crystal. In this case, the model's prediction correctly identifies two particle types, which is more accurate than using one single particle type as the label. Particles such as in (d) are often not straightforward when weathering is incipient (Benet et al., 2024).

juvenile particles, which are critical to avoid overestimating the amount of juvenile component (D’Oriano et al., 2022) and their identification typically requires examination in the SEM (D’Oriano et al., 2014). The high score suggests that the ViT can discriminate between them to some extent (Figure 15e), but a more robust labeling by a team of experts and a larger dataset containing SEM images is necessary to obtain more robust conclusions. On the other hand, the juvenile particles consist of glossy, smoothed surfaces, and vesicular, elongated glass shards and are accurately classified (Figure 15f).

4. The ViT accurately classifies ash particles from plinian and subplinian eruptive activity types (*macro* *F1-score* of 0.95), including free crystals (*F1-score* of 0.92), altered material (*F1-score* of 0.93) and juvenile (1.0), but less accurate for lithics (*F1-score* of 0.77). The juvenile particles consist of fragments of pumice



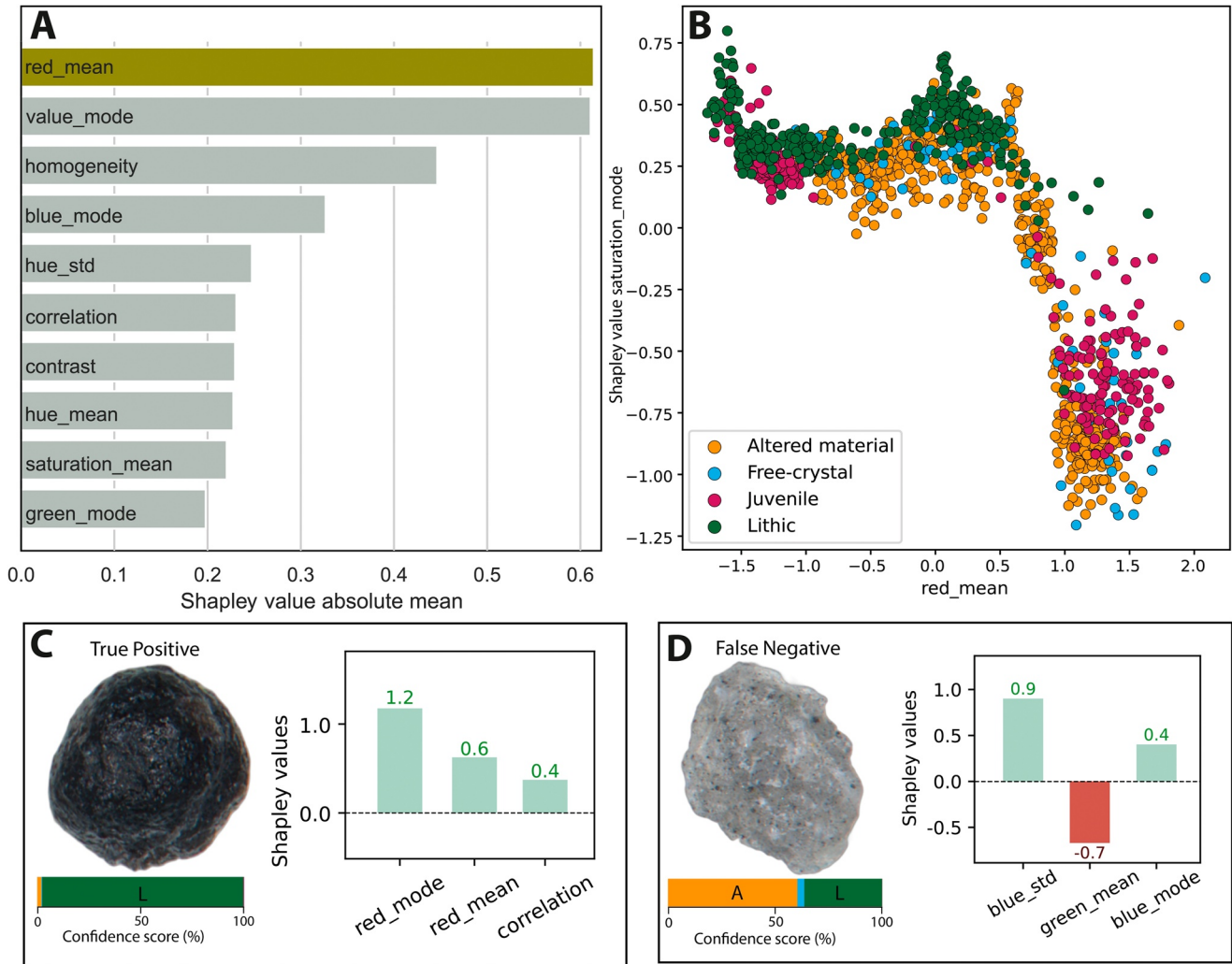
**Figure 10.** Summary plots to illustrate the features that contribute the most to the predictions of the juvenile particles. (a) Feature importance based on the mean of the Shapley values. (b) Shapley dependence plot. Note a cluster of juvenile particles around *saturation\_mode* values between 1 and 3. (c, d) are examples of two predictions of the particle image, with the horizontal bar showing the confidence score across particle types, and the vertical bars the associated Shapley values. (c) Shows a True Positive predicted at maximum confidence score with the *hue\_mean* (chromaticity), *saturation\_mode* (mode of the intensity of the color), and *convexity*. (d) Is an example of a particle that was predicted by the XGBoost model as lithic with a confidence of 70% (size of the green area in horizontal bar plot) based on the *red\_mean* (mean of the red channel), which is predominantly discriminant of lithic particles (a) but was classified as juvenile in VolcAshDB.

and all particles are successfully classified (Figure 15g). In contrast, the lithic particles mostly consist of dull gray fragments with rounded edges, and most of the FN are classified as altered material, which may reflect the challenge of classifying particles with incipient weathering into weathered material or lithic (Figure 15h).

## 4. Discussion

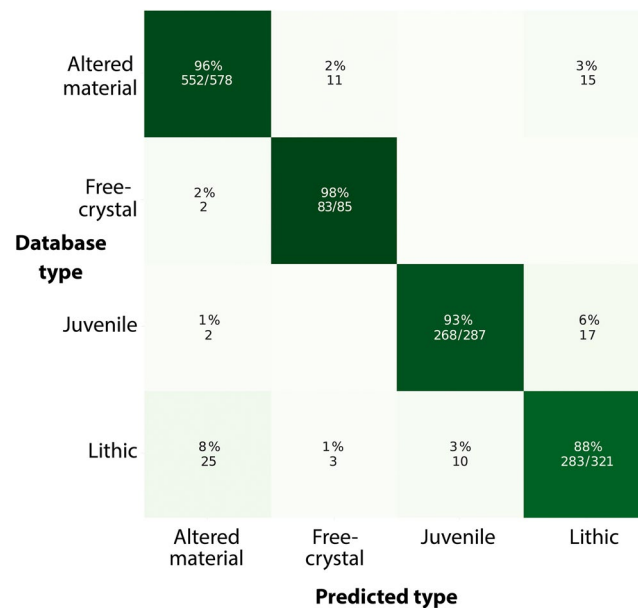
### 4.1. Comparison of Particle Classification Using Features Versus Images

We found that overall the ViT classifies more accurately with particle images (0.93 of *macro F1-score*) than the XGBoost classifies with the particle features (0.77 of *macro F1-score*). This difference is unlikely to be the XGBoost model itself, which is very popular in the literature and has had the best performance in complex



**Figure 11.** Summary plots to explain the predictions of the lithic type. (a) Ranking of the features according to the mean of the Shapley values. (b) The Shapley dependence plot shows correct predictions of lithic particles with high Shapley values at negative values of *value\_mode*. (c, d) show for each prediction: the particle image, the confidence score across particle types, and the associated Shapley values. (c) Shows a dark particle that is correctly classified as lithic with low *value\_mode* (luminosity), whereas (d) shows that XGBoost gives similar confidence scores to the altered material and lithic types, with the former being slightly preferred given the values of *green\_mean*, which are uncharacteristic of the lithic type (shown by negative Shapley value  $-0.7$ ).

classification tasks (Brownlee, 2016; Chen & Guestrin, 2016; Dhaliwal et al., 2018). One possibility is that the extracted features do not retain certain discriminant information from the images, and as a result, the XGBoost is unable to classify particles such as free crystals (0.57 of *F1-score*). On the other hand, maintaining the physical information associated with features makes the model's outcomes more interpretable (e.g., in classification of volcano-seismic signals; Falcin et al., 2021; Malfante et al., 2018) with xAI methods. This is an important advantage over Vision Transformers, whose main xAI tool consists of a heatmap of the region(s) of attention by the model (Dosovitskiy et al., 2020) but appears insufficient to obtain well-founded classification insights for ash particles (Figure 16).



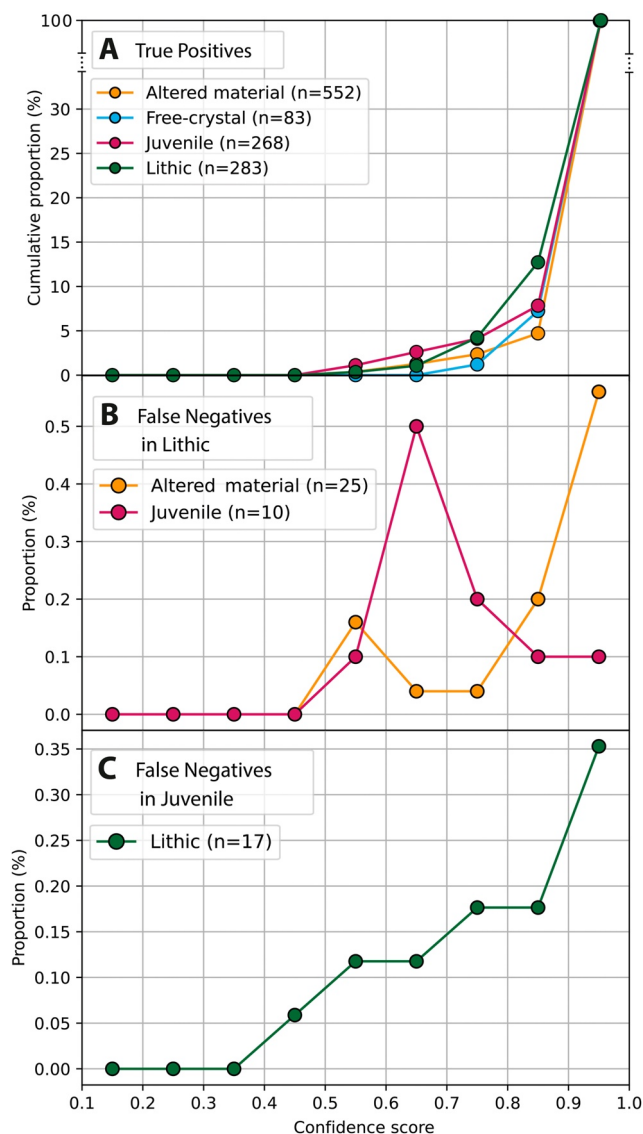
**Figure 12.** Confusion matrix of the predictions by the ViT image classifier. The percentages show the True Positive rate if positioned in the diagonal matrix (darker green), and otherwise, the False-Negative rate (lighter), all percentages with the corresponding number of particles per predicted type. The best classification is for free-crystal followed by altered material, juvenile and lithic.

#### 4.2. Insights From XGBoost to Better Develop the Classification Criteria for Particles Observed With the Binocular Microscope

The XGBoost model gave a medium-to-high classification performance with *macro F1-score* of 0.77, and using the Shapley values, we identified the most discriminant features of each particle type (Table 4). For instance, lithic particles can be distinguished with low values of *value\_mode* which correspond to the luster of the particle. This finding agrees with previous studies that use a dull luster (which corresponds to low values of *value\_mode*) to identify lithic particles (Miwa et al., 2013). On the other hand, juvenile particles have high Shapley values for the *saturation\_mode*. This feature is related to high color intensities as observed under the binocular, but it was not previously recognized as a diagnostic observation of the particle type. These two examples belong to particle types that are well classified and for which the Shapley values are reliable. Shapley values obtained from particles that yielded lower accuracies, such as the free crystals, are not reliable, and thus overall performances should be improved. This could be achieved by enhancing the quality and quantity of VolcAshDB dataset by (a) adding particles to balance the dataset, (b) refining the particle contour in the multi-focused images, so that shape features can measure micro-scaled cavities (Benet et al., 2024), and (c) the inclusion of a new feature that measures the density of lines on the surface, which could be sensitive to planar structures of free crystals.

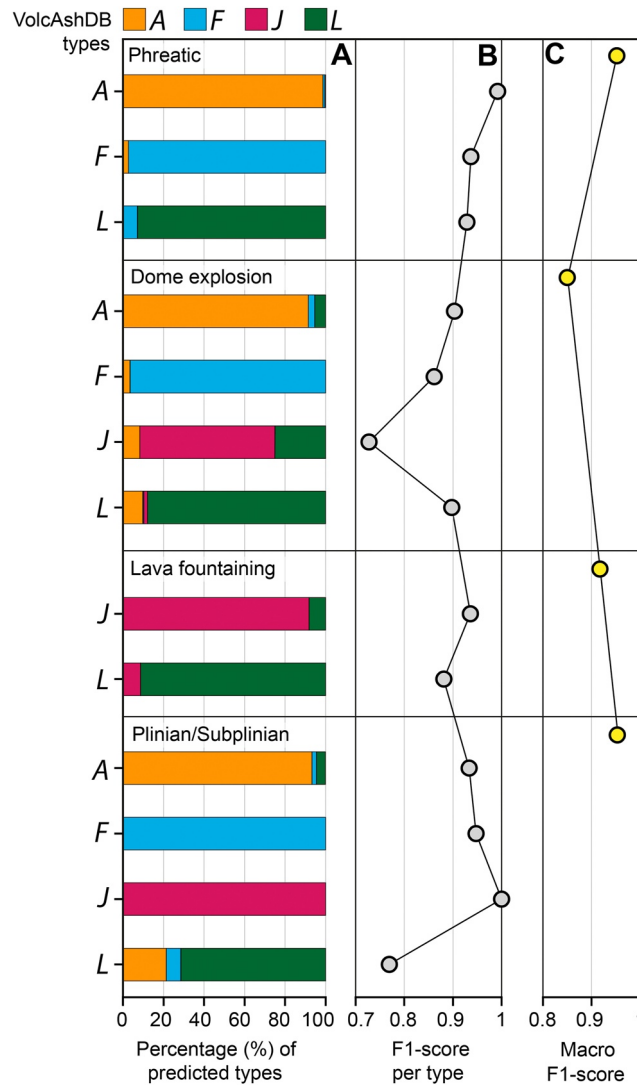
#### 4.3. Deploying the ViT for Automatic Particle Classification

A main goal of our research is to obtain a classifier of ash particles that is as accurate as possible, and which can be applied to objectively classify new datasets in a reproducible manner. The ViT model (*macro F1-score* of 0.93) currently performs very accurately for some samples (e.g., Soufrière de Guadeloupe; *macro F1-score* of 0.95) but is less accurate for others (e.g., Merapi; *macro F1-score* of 0.80). This variation is also found within subgroups of particles. For instance, elongated, highly vesicular, glossy particles from basaltic lava fountaining (Cumbre Vieja, 2021) or pumice fragments (Kelud, 2014) are very accurately



**Figure 13.** Line plots of the confidence score versus (a) the cumulative proportion of True Positives (TP), (b) False Negatives (FN) in free-crystal and (c) lithic types. The distribution of the data has been plotted into 9 bins of size 0.1. We do not use cumulative proportion in (b, c) given the limited number of FNs. Two examples on how to read (a) are described in Figure 6. Note that the ViT predicts True Positives at high confidence score values, although it is less certain about the lithic particle type.

classified, but high crystallinity, blocky, dark particles from dome explosions (Nevados de Chillán, 2016–2018) are not. These changes in classification scores may be due to differences in the particle-forming processes: juvenile particles from Plinian eruptions are originated from a main and short fragmentation episode (Cioni et al., 2015), whereas juvenile particles from dome explosions originate from magma with a long and complex story of slow conduit ascent, degassing, crystallization, fracturing, and recycling (Calder et al., 2015). Moreover, the variability of *F1-scores* between eruptive activity types suggests that to obtain a more robust model that allows for generalization, we need more particles from such problematic subgroups and labeling done by a team of experts. It is also necessary to increase the range of samples, including

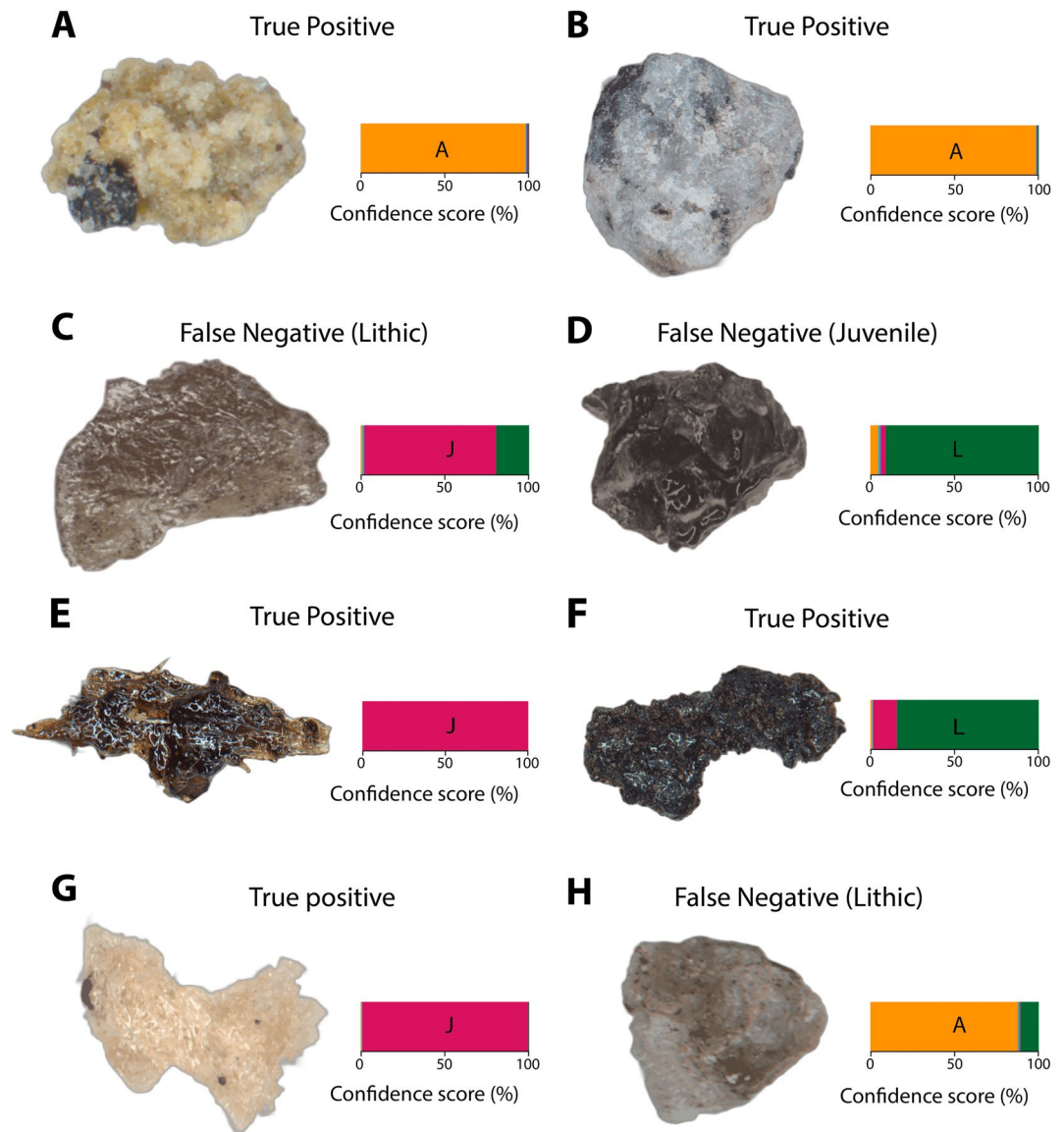


**Figure 14.** (a) Bar charts showing the percentage of predicted types for each particle type in VolcAshDB. If all predictions were the same as in the database, each bar would be single-colored as follows: orange for altered material (*A*), light blue for free-crystal (*F*), magenta for juvenile (*J*), and dark green for lithic (*L*). (b) Shows the *F1-score* for each particle type across eruptive activity types, whereas (c) shows the value of the *macro F1-score* per eruptive activity type. Note the range in *macro F1-score* values (c) from 0.85 for dome explosion to 0.91 for lava fountaining up to 0.95 for phreatic, subplinian and plinian eruptive activity types. The exact values of this figure can be found in Table S11 in Supporting Information S1.

eruptive styles like strombolian and phreatomagmatic activity, and andesitic magma compositions, are probably needed the most.

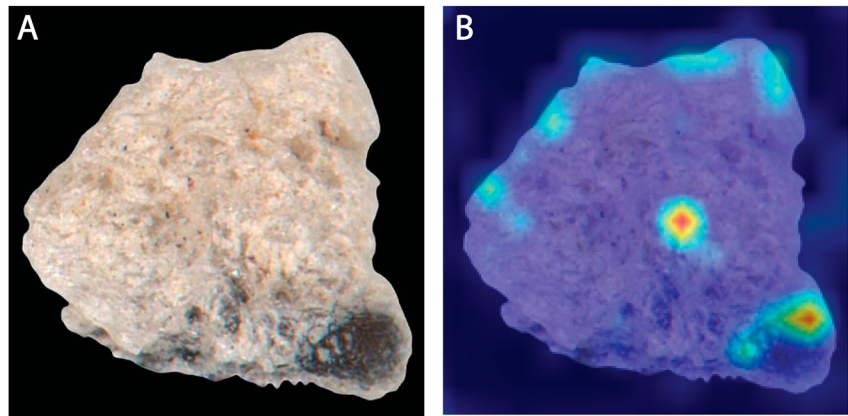
#### 4.4. A ViT Particle Classifier for Volcano Monitoring

Volcano observatories and laboratories are often equipped with binocular microscopes that can acquire standard, single-focus particle images. Our dataset and analysis are based on multi-focused images and therefore it is not clear how directly applicable it is for single-focus images. To evaluate this, we performed a preliminary test of ViT's ability to classify single-focus images from a small dataset of ~1,200 images from Nevados de Chillán

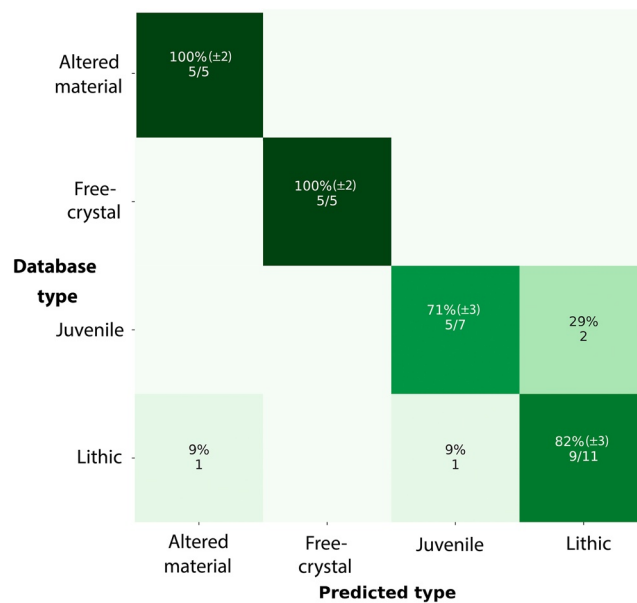


**Figure 15.** Representative examples of particle images and the predictions and their associated confidence score across eruptive activities, including phreatic (a, b), dome explosion (c, d), lava fountaining (e, f), and subplinian/plinian (g, h). Note that False Negatives contain in brackets the particle type according to VolcAshDB, and that the color code and symbol letter are the same as in Figure 14.

(Benet et al., 2021). The dataset comprises 400 particles, with 3 images per particle at different focus depths. Using the same split ratio (80:20) would yield a very small training set, and thus we used all particles for training, except for 28 representative particles as described in Benet et al. (2021) as a test. Fine-tuning the ViT took only 3 hr and we obtained relatively high accuracies (*macro F1-score* of 0.84) on the test set (Figure 17). This suggests that volcano observatories could potentially use a ViT and obtain an objective score on a particle-by-particle basis relatively rapidly.



**Figure 16.** Example of (a) one multi-focused binocular image of a pumice particle from Mount St. Helens, 1980, which is overlain by (b) a heatmap of the regions of attention by the base Vision Transformer (Dosovitskiy et al., 2020), is typically used for interpreting image classifier's predictions. It does not appear easy to discern which aspects of the particle were relevant for classification.



**Figure 17.** Confusion matrix of the predictions by the ViT image classifier after being fine-tuned with a single-focused, small training set (~370 particles from Benet et al. (2021)). The percentages show the True Positive rate if positioned in the diagonal matrix (darker green), and otherwise, the False-Negative rate (lighter), all percentages with the corresponding number of particles per predicted type. Note that we only used 28 particles for the test set (see more details in the main text) and obtained a *macro F1-score* of 0.84. Since the subset is small, we report an error as the square root of the number of particles, which is known in statistics as the implicit random error (Ahmed, 2015).

## 5. Conclusions

Classification of the different particles that make up volcanic ash is often difficult because diagnostic criteria are not standardized and thus reliable, and systematic identification of a given particle type is not straightforward. In this contribution, we attempt to alleviate this situation by exploring the use of state-of-the-art machine learning-based models to identify the most discriminant features of each particle type and to evaluate their ability to classify particles. The identified features provide new insights into the recognition of juvenile and lithic particles toward a standardized classification. The image classifier performs at very high accuracies, although the variability across eruptions and types shows that its capability to generalize to new samples is still unclear. Higher numbers of particles from a wider variety of eruptions and volcanoes into VolcAshDB coupled to ML models should allow for unbiased comparison of ash samples, and reproducible classification of their particles as a tool for volcano monitoring studies.

## Data Availability Statement

Particle images and features can be downloaded through the publicly available VolcAshDB web database at <http://volcashedb.ipgp.fr/>. Details on the feature measurement and image acquisition are described in Benet et al., 2024. The used feature dataset can be found in the GitHub repository [https://github.com/dbenet-max/volcashedb\\_classification](https://github.com/dbenet-max/volcashedb_classification) (dbenet-max & DBenet, 2023) under the name “qia\_processed.csv.” The repository also contains two relevant codes: the Python code for hyperparameter optimization, development, and interpretation via xAI of the XGBoost, and the code for deployment via the API Hugging Face of the ViT.

## References

- Ahmed, S. N. (2015). Essential statistics for data analysis. In *Physics and engineering of radiation detection*. <https://doi.org/10.1016/b978-0-12-801363-2.00009-7>
- Alvarado, G. E., Mele, D., Dellino, P., de Moor, J. M., & Avaró, G. (2016). Are the ashes from the latest eruptions (2010–2016) at Turrialba volcano (Costa Rica) related to phreatic or phreatomagmatic events? *Journal of Volcanology and Geothermal Research*, 327, 407–415. <https://doi.org/10.1016/j.jvolgeores.2016.09.003>
- Ayyadevara, V. K., & Reddy, Y. (2020). *Modern computer vision with PyTorch: Explore deep learning concepts and implement over 50 real-world image applications*. Packt Publishing Ltd.
- Bebbington, M. S., & Jenkins, S. F. (2019). Intra-eruption forecasting. *Bulletin of Volcanology*, 81(6), 1–17. <https://doi.org/10.1007/s00445-019-1294-9>
- Benet, D., Costa, F., Pedreros, G., & Cardona, C. (2021). The volcanic ash record of shallow magma intrusion and dome emplacement at Nevados de Chillán Volcanic complex, Chile. *Journal of Volcanology and Geothermal Research*, 417, 107308. <https://doi.org/10.1016/j.jvolgeores.2021.107308>
- Benet, D., Costa, F., Widiwijayanti, C., Pallister, J., Pedreros, G., Allard, P., et al. (2024). VolcAshDB: A volcanic ash DataBase of classified particle images and features. *Bulletin of Volcanology*, 86(1), 1–30. <https://doi.org/10.1007/s00445-023-01695-4>
- Biass, S., Jenkins, S. F., Aeberhard, W. H., Delmelle, P., & Wilson, T. (2022). Insights into the vulnerability of vegetation to tephra fallouts from interpretable machine learning and big Earth observation data. *Natural Hazards and Earth System Sciences*, 22(9), 2829–2855. <https://doi.org/10.5194/nhess-22-2829-2022>
- Brownlee, J. (2016). XGBoost with python: Gradient boosted trees with XGBoost and scikit-learn. *Machine Learning Mastery*.
- Brownlee, J. (2020). Imbalanced classification with python. *Machine Learning Mastery*, 463.
- Calder, E. S., Lavallée, Y., Kendrick, J. E., & Bernstein, M. (2015). Lava dome eruptions. *The Encyclopedia of Volcanoes*, 343–362. <https://doi.org/10.1016/B978-0-12-385938-9.00018-3>
- Cashman, K. V., & Hoblitt, R. P. (2004). Magmatic precursors to the 18 May 1980 eruption of Mount St. Helens, USA. *Geology*, 32(2), 141–144. <https://doi.org/10.1130/G20078.1>
- Chen, T., & Guestrin, C. (2016). XGBoost: A scalable tree boosting system. In *Proceedings of the ACM SIGKDD international conference on knowledge discovery and data mining* (Vol. 13–17, pp. 785–794). <https://doi.org/10.1145/2939672.2939785>
- Cioni, R., Pistolesi, M., Bertagnini, A., Bonadonna, C., Hoskuldsson, A., & Scateni, B. (2014). Insights into the dynamics and evolution of the 2010 Eyjafjallajökull summit eruption (Iceland) provided by volcanic ash textures. *Earth and Planetary Science Letters*, 394, 111–123. <https://doi.org/10.1016/j.epsl.2014.02.051>
- Cioni, R., Pistolesi, M., & Rosi, M. (2015). Plinian and Subplinian eruptions. In *The encyclopedia of volcanoes* (2nd ed.). Elsevier Inc. <https://doi.org/10.1016/B978-0-12-385938-9.00029-8>
- dbenet-max/DBenet. (2023). dbenet-ntu/volcashedb\_classification: Classification of VolcAshDB with ML (v0.0) [Dataset]. Zenodo. <https://doi.org/10.5281/zenodo.10409409>
- Dellino, P., & La Volpe, L. (1996). Image processing analysis in reconstructing fragmentation and transportation mechanisms of pyroclastic deposits. The case of Monte Pilato-Rocche Rosse eruptions, Lipari (Aeolian Islands, Italy). *Journal of Volcanology and Geothermal Research*, 71(1), 13–29. [https://doi.org/10.1016/0377-0273\(95\)00062-3](https://doi.org/10.1016/0377-0273(95)00062-3)
- Deng, J., Dong, W., Socher, R., Li, L., Li, K., & Fei-Fei, L. (2009). ImageNet: A large-scale hierarchical image database (pp. 248–255). <https://doi.org/10.1109/cvprw.2009.5206848>
- Dhaliwal, S. S., Nahid, A. A., & Abbas, R. (2018). Effective intrusion detection system using XGBoost. *Information*, 9(7), 149. <https://doi.org/10.3390/info9070149>

## Acknowledgments

D.B. is grateful to C. Bouvet de Maisonneuve, J. Pallister, J. Taddeucci and A. Rust for insightful discussions, Do Xuan Long for help on the use of the Hugging Face API for image classification, and S. Biass for advice and help on the use of the SHAP method for this study, and to E. Tan for support on the Gekko cluster. The manuscript was revised by three anonymous reviewers and editor M. Edmonds. This research was supported by the Earth Observatory of Singapore via funding from the National Research Foundation Singapore and the Singapore Ministry of Education under the Research Centres of Excellence initiative. F.C. acknowledges also support by a Chaire d'Excellence grant by the Université Paris Cité.

- D'Oriano, C., Bertagnini, A., Cioni, R., & Pompilio, M. (2014). Identifying recycled ash in basaltic eruptions. *Scientific Reports*, *4*(1), 5851. <https://doi.org/10.1038/srep05851>
- D'Oriano, C., Del Carlo, P., Andronico, D., Cioni, R., Gabellini, P., Cristaldi, A., & Pompilio, M. (2022). Syn-eruptive processes during the January–February 2019 ash-rich emissions cycle at Mt. Etna (Italy): Implications for petrological monitoring of volcanic ash. *Frontiers in Earth Science*, *10*. <https://doi.org/10.3389/feart.2022.824872>
- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., et al. (2020). An image is worth 16×16 Words: Transformers for image recognition at scale.
- Duchi, J., Hazan, E., & Singer, Y. (2011). Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, *12*, 2121–2159.
- Dürig, T., Bowman, M. H., White, J. D. L., Murch, A., Mele, D., Verolino, A., & Dellino, P. (2018). Particle shape analyzer Partisan—An open source tool for multi-standard two-dimensional particle morphometry analysis. *Annals of Geophysics*, *61*(6). <https://doi.org/10.4401/ag-7865>
- Dürig, T., Ross, P. S., Dellino, P., White, J. D. L., Mele, D., & Comida, P. P. (2021). A review of statistical tools for morphometric analysis of juvenile pyroclasts. *Bulletin of Volcanology*, *83*(11), 79. <https://doi.org/10.1007/s00445-021-01500-0>
- Falcin, A., Métaixian, J. P., Mars, J., Stutzmann, É., Komorowski, J. C., Moretti, R., et al. (2021). A machine-learning approach for automatic classification of volcanic seismicity at La Soufrière Volcano, Guadeloupe. *Journal of Volcanology and Geothermal Research*, *411*, 107151. <https://doi.org/10.1016/j.jvolgeores.2020.107151>
- Feuillard, M., Allegre, C. J., Brandeis, G., Gaulon, R., Le Mouel, J., Mercier, J. C., et al. (1983). The 1975–1977 crisis of la Soufriere de Guadeloupe (F.W.I): A still-born magmatic eruption. *Journal of Volcanology and Geothermal Research*, *16*(3–4), 317–334. [https://doi.org/10.1016/0377-0273\(83\)90036-7](https://doi.org/10.1016/0377-0273(83)90036-7)
- Fürnkranz, J., Hüllermeier, E., Loza Mencía, E., & Brinker, K. (2008). Multilabel classification via calibrated label ranking. *Machine Learning*, *73*(2), 133–153. <https://doi.org/10.1007/s10994-008-5064-8>
- Gaunt, H. E., Bernard, B., Hidalgo, S., Proaño, A., Wright, H., Mothes, P., et al. (2016). Juvenile magma recognition and eruptive dynamics inferred from the analysis of ash time series: The 2015 reawakening of Cotopaxi volcano. *Journal of Volcanology and Geothermal Research*, *328*, 134–146. <https://doi.org/10.1016/j.jvolgeores.2016.10.013>
- Géron, A. (2017). Hands-on machine learning with Scikit-Learn and TensorFlow: Concepts, tools, and techniques to build intelligent systems. In *Hands-on machine learning with Scikit-Learn and TensorFlow: Concepts, tools, and techniques to build intelligent systems*.
- Gianfagna, L., & Di Cecco, A. (2021). *Explainable AI with python* (pp. 1–202). Springer.
- Hall-Beyer, M. (2017). GLCM texture: A tutorial. *17th International Symposium on Ballistics*, *2*(March), 18–19.
- Haralick, R. M., Dinstein, I., & Shanmugam, K. (1973). Textural features for image classification. In *IEEE transactions on systems, man and cybernetics, SMC-3* (pp. 610–621). <https://doi.org/10.1109/TSMC.1973.4309314>
- He, K., Zhang, X., Ren, S., & Sun, J. (2015). Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE international conference on computer vision* (pp. 1026–1034).
- He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition. In *Proceedings of the IEEE computer society conference on computer vision and pattern recognition* (pp. 770–778). <https://doi.org/10.1109/CVPR.2016.90>
- Herrera, F., Chartre, F., Rivera, A. J., & del Jesus, M. J. (2016). *Multi-label classification*. Springer. <https://doi.org/10.4018/jdwm.2007070101>
- Hincks, T. K., Komorowski, J. C., Sparks, S. R., & Aspinall, W. P. (2014). Retrospective analysis of uncertain eruption precursors at La Soufrière volcano, Guadeloupe, 1975–1977: Volcanic hazard assessment using a Bayesian Belief Network approach. *Journal of Applied Volcanology*, *3*(1), 3. <https://doi.org/10.1186/2191-5040-3-3>
- Hughes, I., & Hase, T. (2010). *Measurements and their uncertainties: A practical guide to modern error analysis*. OUP Oxford.
- Kandel, I., & Castelli, M. (2020). The effect of batch size on the generalizability of the convolutional neural networks on a histopathology dataset. *ICT Express*, *6*(4), 312–315. <https://doi.org/10.1016/j.ict.2020.04.010>
- Kondylatos, S., Prapas, I., Ronco, M., Papoutsis, I., Camps-Valls, G., Piles, M., et al. (2022). Wildfire danger prediction and understanding with deep learning. *Geophysical Research Letters*, *49*(17), 1–11. <https://doi.org/10.1029/2022GL099368>
- Kotsiantis, S. B. (2013). Decision trees: A recent overview. *Artificial Intelligence Review*, *39*(4), 261–283. <https://doi.org/10.1007/s10462-011-9272-4>
- Lee, J. J., Aime, M. C., Rajwa, B., & Bae, E. (2022). Machine learning-based classification of mushrooms using a smartphone application. *Applied Sciences*, *12*(22), 11685. <https://doi.org/10.3390/app122211685>
- Le Guern, F., Bernard, A., & Chevrier, R. M. (1980). Soufrière de Guadeloupe 1976–1977 eruption—Mass and energy transfer and volcanic health hazards. *Bulletin Volcanologique*, *43*(3), 577–593. <https://doi.org/10.1007/BF02597694>
- Leibrandt, S., & Le Pennec, J. L. (2015). Towards fast and routine analyses of volcanic ash morphometry for eruption surveillance applications. *Journal of Volcanology and Geothermal Research*, *297*, 11–27. <https://doi.org/10.1016/j.jvolgeores.2015.03.014>
- Liu, E. J., Cashman, K. V., & Rust, A. C. (2015). Optimising shape analysis to quantify volcanic ash morphology. *GeoResJ*, *8*, 14–30. <https://doi.org/10.1016/j.grj.2015.09.001>
- Liu, Z., Mao, H., Wu, C.-Y., Feichtenhofer, C., Darrell, T., & Xie, S. (2022). A ConvNet for the 2020s (pp. 11966–11976). <https://doi.org/10.1109/cvpr52688.2022.01167>
- Loshchilov, I., & Hutter, F. (2019). Decoupled weight decay regularization. In *7th international conference on learning representations, ICLR 2019*.
- Lundberg, S. M., Erion, G. G., & Lee, S.-I. (2018). Consistent individualized feature attribution for tree ensembles (Vol. 2). Retrieved from <http://arxiv.org/abs/1802.03888>
- Lundberg, S. M., & Lee, S.-I. (2017). A unified approach to interpreting model predictions. *Advances in Neural Information Processing Systems*, *30*.
- Maeno, F., Nakada, S., Yoshimoto, M., Shimano, T., Hokanishi, N., Zaennudin, A., & Iguchi, M. (2019). A sequence of a plinian eruption preceded by dome destruction at Kelud volcano, Indonesia, on February 13, 2014, revealed from tephra fallout and pyroclastic density current deposits. *Journal of Volcanology and Geothermal Research*, *382*, 24–41. <https://doi.org/10.1016/j.jvolgeores.2017.03.002>
- Malfante, M., Dalla Mura, M., Métaixian, J. P., Mars, J. I., Macedo, O., & Inza, A. (2018). Machine learning for volcano-seismic signals: Challenges and perspectives. *IEEE Signal Processing Magazine*, *35*(2), 20–30. <https://doi.org/10.1109/msp.2017.2779166>
- Mandal, S., Mones, S. M. B., Das, A., Balas, V. E., Shaw, R. N., & Ghosh, A. (2021). Single shot detection for detecting real-time flying objects for unmanned aerial vehicle. In *Artificial intelligence for future generation robotics*. INC. <https://doi.org/10.1016/B978-0-323-85498-6.00005-8>
- Marzocchi, W., Newhall, C., & Woo, G. (2012). The scientific management of volcanic crises. *Journal of Volcanology and Geothermal Research*, *247–248*, 181–189. <https://doi.org/10.1016/j.jvolgeores.2012.08.016>
- Mishra, P. (2022). Practical explainable AI using python. In *Practical explainable AI using python*. <https://doi.org/10.1007/978-1-4842-7158-2>

- Miwa, T., Geshi, N., & Shinohara, H. (2013). Temporal variation in volcanic ash texture during a vulcanian eruption at the sakurajima volcano, Japan. *Journal of Volcanology and Geothermal Research*, 260, 80–89. <https://doi.org/10.1016/j.jvolgeores.2013.05.010>
- Miwa, T., Toramaru, A., & Iguchi, M. (2009). Correlations of volcanic ash texture with explosion earthquakes from vulcanian eruptions at Sakurajima volcano, Japan. *Journal of Volcanology and Geothermal Research*, 184(3–4), 473–486. <https://doi.org/10.1016/j.jvolgeores.2009.05.012>
- Miyagi, I., Geshi, N., Hamasaki, S., Oikawa, T., & Tomiya, A. (2020). Heat source of the 2014 phreatic eruption of Mount Ontake, Japan. *Bulletin of Volcanology*, 82(4), 33. <https://doi.org/10.1007/s00445-020-1358-x>
- Molnar, C., Li, J., Kim, J. S., Plumb, G., & Talwalkar, A. (2021). Interpretable machine learning. *ACM Queue*, 19(6), 28–56. <https://doi.org/10.1145/3511299>
- Moran, S. C., Newhall, C., & Roman, D. C. (2011). Failed magmatic eruptions: Late-stage cessation of magma ascent. *Bulletin of Volcanology*, 73(2), 115–122. <https://doi.org/10.1007/s00445-010-0444-x>
- Newhall, C. G., & Punongbayan, R. S. (1996). The narrow margin of successful volcanic-risk mitigation. In *Monitoring and mitigation of volcano hazards* (pp. 807–838). Springer Science and Business Media.
- Nurfiani, D., & Bouvet de Maisonneuve, C. (2018). Furthering the investigation of eruption styles through quantitative shape analyses of volcanic ash particles. *Journal of Volcanology and Geothermal Research*, 354, 102–114. <https://doi.org/10.1016/j.jvolgeores.2017.12.001>
- Owen, L. (2022). *Hyperparameter tuning with Python: Boost your machine learning model's performance via hyperparameter tuning*. Packt Publishing Ltd.
- Paladio-Melasantos, M. L., Solidum, R. U., Scott, W. E., Quiambao, R. B., Umbal, J. V., Rodolfo, K. S., et al. (1996). Tephra falls of the 1991 eruptions of Mount Pinatubo. In C. G. Newhall (Ed.), (Eds), *Others, fire and mud; eruptions and Lahars of Mount Pinatubo, Philippines*, Philippine Institute of Volcanology and Seismology, Quezon City, layer D (pp. 413–535). <https://doi.org/10.1159/000153100>
- Panati, C., Wagner, S., & Bruggenwirth, S. (2022). Feature relevance evaluation using grad-CAM, LIME and SHAP for deep learning SAR data classification. *Proceedings International Radar Symposium*, 457–462.
- Pardo, N., Avellaneda, J. D., Rausch, J., Jaramillo-Vogel, D., Gutiérrez, M., & Foubert, A. (2020). Decrypting silicic magma/plug fragmentation at Azufral Crater Lake, Northern Andes: Insights from fine to extremely fine ash morpho-chemistry. *Bulletin of Volcanology*, 82(12), 79. <https://doi.org/10.1007/s00445-020-01418-z>
- Pardo, N., Cronin, S. J., Németh, K., Brenna, M., Schipper, C. I., Breard, E., et al. (2014). Perils in distinguishing phreatic from phreatomagmatic ash; insights into the eruption mechanisms of the 6 August 2012 Mt. Tongariro eruption, New Zealand. *Journal of Volcanology and Geothermal Research*, 286, 397–414. <https://doi.org/10.1016/j.jvolgeores.2014.05.001>
- Re, G., Corsaro, R. A., D'Oriano, C., & Pompilio, M. (2021). Petrological monitoring of active volcanoes: A review of existing procedures to achieve best practices and operative protocols during eruptions. *Journal of Volcanology and Geothermal Research*, 419, 107365. <https://doi.org/10.1016/j.jvolgeores.2021.107365>
- Romero, J. E., Burton, M., Cáceres, F., Taddeucci, J., Civico, R., Ricci, T., et al. (2022). The initial phase of the 2021 Cumbre Vieja ridge eruption (Canary Islands): Products and dynamics controlling edifice growth and collapse. *Journal of Volcanology and Geothermal Research*, 431(July), 107642. <https://doi.org/10.1016/j.jvolgeores.2022.107642>
- Ross, P. S., Dürig, T., Comida, P. P., Lefebvre, N., White, J. D. L., Andronico, D., et al. (2022). Standardized analysis of juvenile pyroclasts in comparative studies of primary magma fragmentation; 1. Overview and workflow. *Bulletin of Volcanology*, 84(1), 1–29. <https://doi.org/10.1007/s00445-021-01516-6>
- Rowe, M. C., Thornber, C. R., & Kent, A. J. R. (2008). Identification and evolution of the juvenile component in. In *A volcano rekindled: The renewed eruption of Mount St. Helens, 2004–2006* (pp. 2004–2006).
- Scheidegger, K. F., Federman, A. N., & Tallman, A. M. (1982). Compositional heterogeneity of tephra from the 1980 eruptions of Mount St. Helens. *Journal of Geophysical Research*, 87(B13), 10861–10881. <https://doi.org/10.1029/jb087ib13p10861>
- Shapley, L. S. (1953). A value for n-person games. In *Contributions to the theory of games II* (pp. 307–318). <https://doi.org/10.1515/9781400881970-018>
- Shoji, D., Noguchi, R., Otsuki, S., & Hino, H. (2018). Classification of volcanic ash particles using a convolutional neural network and probability. *Scientific Reports*, 8(1), 8111. <https://doi.org/10.1038/s41598-018-26200-2>
- Sujatha, R., Chatterjee, J. M., Jhanjhi, N. Z., & Brohi, S. N. (2021). Performance of deep learning vs machine learning in plant leaf disease detection. *Microprocessors and Microsystems*, 80, 103615. <https://doi.org/10.1016/j.micpro.2020.103615>
- Sutskever, I., Martens, J., Dahl, G., & Hinton, G. (2013). On the importance of initialization and momentum in deep learning. In *30th international conference on machine learning, ICML 2013, PART 3* (pp. 2176–2184).
- Suzuki, Y., Nagai, M., Maeno, F., Yasuda, A., Hokanishi, N., Shimano, T., et al. (2013). Precursory activity and evolution of the 2011 eruption of Shinmoe-dake in Kirishima volcano—insights from ash samples. *Earth Planets and Space*, 65(6), 591–607. <https://doi.org/10.5047/eps.2013.02.004>
- Taddeucci, J., Pompilio, M., & Scarlato, P. (2002). Monitoring the explosive activity of the July–August 2001 eruption of Mt. Etna (Italy) by ash characterization. *Geophysical Research Letters*, 29(8), 1–4. <https://doi.org/10.1029/2001GL014372>
- Tang, S., Ghorbani, A., Yamashita, R., Rehman, S., Dunnmon, J. A., Zou, J., & Rubin, D. L. (2021). Data valuation for medical imaging using Shapley value and application to a large-scale chest X-ray dataset. *Scientific Reports*, 11(1), 1–9. <https://doi.org/10.1038/s41598-021-87762-2>
- Tilling, R. ~I. (2008). The critical role of volcano monitoring in risk reduction. *Advances in Geosciences*, 14, 3–11. <https://doi.org/10.5194/adgeo-14-3-2008>
- Utami, S. B., Costa, F., Lesage, P. H., Allard, P., & Humaida, H. (2021). Fluid fluxing and accumulation drive decadal and short-lived explosive basaltic andesite eruptions preceded by limited volcanic unrest. *Journal of Petrology*, 62(11), 1–29. <https://doi.org/10.1093/ptrology/legab086>
- Verdhan, V. (2020). Supervised learning with python. *Okänd. Irland: Apress*. <https://doi.org/10.1007/978-1-4842-6156-9>
- Watanabe, K., Danhara, T., Watanabe, K., Terai, K., & Yamashita, T. (1999). Juvenile volcanic glass erupted before the appearance of the 1991 lava dome, Unzen volcano, Kyushu, Japan. *Journal of Volcanology and Geothermal Research*, 89(1–4), 113–121. [https://doi.org/10.1016/S0377-0273\(98\)00127-9](https://doi.org/10.1016/S0377-0273(98)00127-9)