



# Automatic classification and segmentation of Snow Micro Penetrometer profiles with machine learning algorithms

Julia Kaltenborn<sup>1,2,3,4</sup>, Amy R. Macfarlane<sup>1</sup>, Viviane Clay<sup>2,5</sup>, and Martin Schneebeli<sup>1</sup>

<sup>1</sup>WSL Institute for Snow and Avalanche Research SLF, Flüelastrasse 11, 7260 Davos Dorf, Switzerland

<sup>2</sup>Institute of Cognitive Science, University Osnabrück, Wachsbleiche 27, 49090 Osnabrück, Germany

<sup>3</sup>Mila - Quebec AI Institute, 6666 Rue Saint-Urbain, QC H2S 3H1, Montréal, Canada

<sup>4</sup>School of Computer Science, McGill University, 3480 Rue University, QC H3A 2A7, Montréal, Canada

<sup>5</sup>Numenta, 889 Winslow Street, CA 94063, Redwood City, United States

**Correspondence:** Julia Kaltenborn (julia.kaltenborn[at]mail.mcgill.ca)

**Abstract.** Snow-layer segmentation and classification is an essential diagnostic task for a wide variety of cryospheric applications. The SnowMicroPen (SMP) measures the snowpack's penetration force at submillimetre resolution against the snow depth. The resulting depth-force profile can be parameterized for density and specific surface area. However, no information on traditional snow types is currently extracted automatically. The labeling of snow types is a time-intensive task that requires practice and becomes infeasible for large datasets. Previous work showed that automated segmentation and classification is in theory possible, but can either not be applied to data straight from the field or needs additional time-costly information, such as from classified snow pits. To address this gap, we evaluate how well machine learning models can automatically segment and classify SMP profiles. We trained fourteen different models, among them semi-supervised models and artificial neural networks (ANNs), on the MOSAiC SMP dataset, a large collection of snow profiles on Arctic sea ice. We found that SMP profiles can be successfully segmented and classified into snow classes, based solely on the SMP's signal. The model comparison provided in this study enables practitioners to choose a model that is suitable for their task and dataset. The findings presented will facilitate and accelerate snow type identification through SMP profiles. Overall, snowdragon creates a link between traditional snow classification and high-resolution force-depth profiles. With such a tool, traditional snow profile observations can be compared to SMP profiles.

## 1 Introduction

The cryosphere covers around 10 % percent of our earth and plays a significant role in stabilizing earth's climate (Pörtner et al., 2019). Snow cover plays a role in optics, heat, and mass balance and is one of the largest uncertainties in global climate models (Sturm and Massom, 2017; Steger et al., 2013; Douville et al., 1995). Snow layer segmentation and classification put forth knowledge about the atmospheric conditions a snowpack has experienced (Colbeck, 1987; Fierz et al., 2009). This knowledge helps to discern fundamental snow and climate mechanisms in the Arctic and to analyze polar tipping points. Classification of snow types is essential to assess the state of our cryosphere and is thus of interest for polar, cryospheric, and climate change research.



Traditionally, snow stratigraphy measurements are made in snow pits. These pits are dug manually, vertically into snowpacks and require trained operators and a substantial time commitment. To accelerate these measurements, the SnowMicroPen (SMP),  
25 a portable high-resolution snow penetrometer, can be used (Johnson and Schneebeli, 1998). Schneebeli and Johnson (1998) have demonstrated the SMP as a capable tool for rapid snow type classification and layer segmentation. The measurement results are stored in an SMP profile that consists of the penetration force signal of the measurement tip in Newton and the depth signal indicating how far the tip moved. Afterwards, the SMP profiles must be manually labeled, which requires time, practice, and becomes infeasible for large datasets. Machine learning (ML) algorithms could be used to automate this process. As a  
30 consequence this would (1) immensely accelerate the SMP analysis, (2) enable the analysis of large datasets, and (3) make the training of interdisciplinary scientists in snow type categorization obsolete.

The nearest neighbor method of Satyawali et al. (2009) was the first model that automated both segmentation and classification of SMP profiles without being dependent on snowpit information. Their algorithm could predict five different snow types, however, their testing dataset was too small to be representative and they excluded data points with uncertain snow  
35 classes. Furthermore, Satyawali et al. (2009) achieved only a high classification performance by including knowledge-based rules which do not generalize on datasets from other regions or seasons.

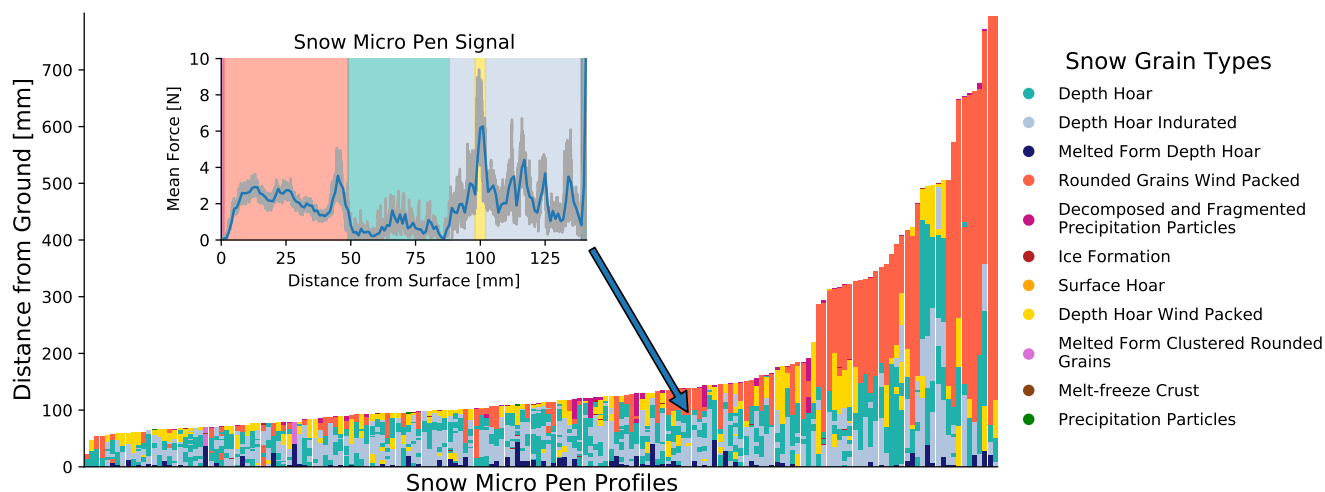
Havens et al. (2012) used previously segmented SMP profiles and classified the snow type of each layer with the help of a random forest model. Their work builds upon their previous work with single decision trees (Havens et al., 2010). Their model could be improved further by adding more than three snow types, allowing also for layers thinner than 100 mm and  
40 most importantly, by automating the segmentation step as well.

The support vector machine (SVM) approach by King et al. (2020) automated both the segmentation and classification, and achieved good accuracy scores for three different snow types. However, they are relying on additional snowpit information to achieve these results.

While all these works put forward the task of automated SMP analysis, SMP practitioners still lack a method that can be  
45 used in practice. Practitioners need a model that fully automates their SMP analysis: (1) without the need of digging a snow pit, (2) picking layers manually or (3) constructing specific knowledge rules. Furthermore, SMP practitioners need models that can deal with SMP profiles coming straight from the field. This implies that (4) the profiles have multiple snow types (more than three) and that (5) no layers are excluded. The aim of this study is to provide models that fully automate SMP analysis and can directly be used in the field, addressing all five mentioned needs.

50 To this end, we implemented fourteen different machine learning (ML) models and compared their performance on the MOSAiC SMP dataset, consisting of 164 labeled profiles (see Fig. 1). Thereby, we provide the first comparable performance overview of different models classifying and segmenting SMP profiles. Moreover, we used – to the best of our knowledge – for the first time semi-supervised methods and artificial neural networks (ANNs) for SMP classification and segmentation.

Results show that especially the artificial neural networks (ANNs), such as the long short-term memory (LSTM) and the  
55 Encoder-Decoder, can produce predictions that are similar to profiles labeled by experts and achieve the best results among all models. However, the choice of the model depends mostly on the individual needs of an SMP user because factors such as explainability, desired sensitivity to rare classes, available time, and computational resources must be taken into consideration.



**Figure 1.** All 164 labeled SnowMicroPen (SMP) profiles used for training, validation (80%), and testing (20%). Each bar represents one SMP profile. The colors encode the different snow grain types. The top of each bar is the air-snow interface and the bottom of each bar is the snow-ground interface of the profile. The in-picture figure illustrates the force signal (grey) and mean force signal (blue) of a single SMP profile (S31H0368). The snow-air interface is on the left, and the bottom of the profile is on the right. The background shading represents the ground-truth labeling of the profile.

The main contributions of this study are:

- Demonstration that SMP profiles straight from the field can be automatically segmented and classified; without manual preparation of the profiles or additional snow-pit data
- Evaluation of semi-supervised models and ANNs for SMP classification and segmentation
- Detailed comparison of different ML models for SMP classification and segmentation

In the following section (Sect. 2) the data and the classification task are described, as well as the fourteen different models that were used in this study. In Sect. 3, the models' performances are presented. Subsequently, the results, their limitations, and future work are discussed in section 4. The impact of this work is addressed in the conclusion (Sect. 5). The code and data availability is outlined directly after the conclusion.

## 2 Methods

### 2.1 Data

All experiments throughout this study used SnowMicroPen profiles from snow on Arctic sea ice. 3680 profiles were collected during the MOSAiC expedition between October 2019 and September 2020 (Nicolaus et al., 2022). 164 profiles from the



cold season (January – May 2020) were labeled and evaluated for this study (see Fig. 1). This study focuses only on profiles of cold snow, as there exists no standardized interpretation of SMP force profiles for wet snow. All profiles collected in the cold season are referred to as “MOSAiC winter data” in the following. The labels indicate which snow type is found at the respective position of the profile. Refer to Fierz et al. (2009) for descriptions of the different snow types referenced here and a classification guideline for snow particles visually observed. The labeling was conducted by a snow expert and is solely based on the properties of the force signal (magnitude, frequency, and gradient) and the signature of the SMP-signal (Schneebeli et al., 1999). After one labeling phase, all profiles were revisited by the same expert to ensure consistent and correct labeling. The surface and the ground of the profiles were detected automatically by the pyngui application of the snowmicropyn package<sup>1</sup>. The labeled profiles were used during training, testing, and validation, while some of the unlabeled profiles were used for semi-supervised models and during generalization tests.

We preprocessed each SMP profile as well as the complete labeled dataset. For each SMP profile we replaced negative force values with 0, summarized the signal into bins (1 mm), and added additional features. During binning we determined mean, variance, maximum, and minimum force signal values. When adding features, time-dependent and location-dependent information is especially relevant: 4 mm and 12 mm sliding windows were applied to extract additional time-dependent information, including variables from the Poisson shot noise model from Löwe and Van Herwijnen (2012). For location-dependent information, we included distance from the ground and the position within the snowpack. We preprocessed the complete labeled dataset by normalizing it, removing profiles from the melting season, and merging snow classes. For example, “Decomposed and Fragmented Precipitation Particles” are merged with the class “Precipitation Particles” since they represent a similar type of snow. The few occurring “Ice Formations” and “Surface Hoar” instances are summarized in the class “Rare”. The data preprocessing ensures that the dataset is clean and that all necessary information, such as time-dependent information, is available during classification.

The resulting dataset has the following properties: (1) There are multiple, noisy, and overlapping classes. (2) There is a between-class imbalance. (3) There is a within-class imbalance, i.e. sub-groups within one class are imbalanced. (4) The labeling of classes is afflicted with uncertainty, i.e. snow experts themselves are not sure to which class exactly some data points belong. The complexity of the data set complicates classification and lowers the maximum achievable accuracy.

## 2.2 Task description

We compare the capabilities of different models to classify and segment the profiles of the MOSAiC winter SMP dataset. To this end, the models first classify each data point of the signal and then summarize the classified points into distinct snow layers (“first-classify-then-segment”). This task can be solved with different learning and classification techniques.

The task can be addressed via **independent classification** or **sequence labeling**. In independent classification, each individual point is classified independently, without looking at other data points. The underlying assumption is that each individual data point carries enough information to be classified solely on that basis. In contrast, sequence labeling assumes that the data

<sup>1</sup><https://snowmicropyn.readthedocs.io/en/latest/>



is an intra-dependent sequence, where the label of each data point also depends on the preceding labels (Nguyen and Guo, 2007).

105 The models can follow either the **supervised**, **unsupervised**, or **semi-supervised learning** regime. In supervised learning, labels are provided to learn an input-output mapping function (Russell and Norvig, 2002). In unsupervised learning, patterns and structure are found in unlabeled data (Ghahramani, 2004), however, no classification is possible, which is why no unsupervised models are employed here. Instead, semi-supervised models are used, which are able to find structures in sparsely labeled data and leverage this information during classification. In the following, all models employed in this work are shortly  
110 presented and put in the context of their learning and task type.

### 2.3 Models

The **majority vote** classifier is used as the baseline for the performance comparison and simply predicts always the majority class (“Rounded Grains Wind Packed”). It satisfies the criteria that a baseline should not require much expertise, should be easy to build, and fast to evaluate (Li et al., 2020).

115 The **cluster-then-predict models** employed in this study, can be separated into three different semi-supervised and independent classification models. Unsupervised methods are used to find clusters in the dataset and subsequently, a supervised model is used to assign labels to the cluster (Soni and Mathai, 2015; Trivedi et al., 2015). As unsupervised model, k-means clustering (Forgy, 1965; Lloyd, 1982), mixture model clustering (GMM) (Bishop, 2006) and Bayesian Gaussian mixture models (BGMM) (Bishop, 2006) were used. The supervised part of the model is a simple majority vote within the clusters, in order  
120 to see if the unsupervised model adds enough information to beat the majority vote baseline.

**Label propagation** is a graph-based, semi-supervised, independent classification algorithm. It propagates the labels of labeled data points to unlabeled ones (Zhu and Ghahramani, 2002). Here, a modified version of this algorithm by Zhou et al. (2004) is used (also known as “label spreading”) (Bengio et al., 2006; Pedregosa et al., 2011).

**Self-trained classifiers** turn a given supervised classifier into a semi-supervised independent classifier. It follows an iterative  
125 approach of training a supervised model on labeled data, predicting more data with the model, and retraining the model with the most confident predictions (Yarowsky, 1995).

**Random forests** (RFs) are ensembles of diversified decision trees (supervised and independent classification). The diversification happens via tree and feature bagging, where only subsets of data or features are used during training (Ho, 1995; Breiman, 2001). Decision trees are simple to build, explainable, white-box classifiers and for these reasons among the most popular machine learning algorithms (Wu et al., 2008). Additionally, a balanced random forest was used with random under-sampling to  
130 balance the data (Chen et al., 2004).

**Support vector machines** (SVMs) construct a hyperplane in a high-dimensional space to solve binary classification tasks (Cortes and Vapnik, 1995; Han et al., 2012) (supervised and independently). When a problem is non-linearly separable, the input data can be projected into a higher-dimensional space until the problem becomes linearly separable. The kernel trick  
135 can be used to circumvent the computationally expensive data transformation involved here. It directly extracts a non-linear optimal hyperplane (Schölkopf et al., 2002).



**K-nearest neighbours** (KNN) is a local, non-parametric classification method that compares samples and classifies new samples based on their  $k$  nearest training data points (supervised and independently). The class of the prediction sample is determined via majority vote. (Fix and Hodges Jr, 1952; Cover and Hart, 1967)

140 **Easy ensemble classifiers** are ensembles of balanced adaptive boosting classifiers (supervised and independent). The method is especially helpful for imbalanced datasets since the learners are trained on different bootstrap samples, which are balanced via random under-sampling. (Liu et al., 2008)

**Long short-term memories** (LSTMs) are a form of artificial neural networks (ANNs) and can perform supervised sequence labeling tasks. ANNs incrementally update their decision function that describes the decision boundary between classes. ANNs  
145 have different nodes, which can be seen as representing different parts of the functions which are weighted differently. During training, the weights of the ANN are optimized by minimizing a loss function via gradient descent. A long short-term memory can handle time-series data. It consists of different memory cells so the LSTM can forget information that is no longer needed, remember information that is required for future decisions, and retrieve information that is required for current decisions. (Hochreiter and Schmidhuber, 1997; Jurafsky and Martin, 2021)

150 **Bidirectional LSTMs** (BLSTMs) connect two independent LSTMs where the first LSTM processes the inputs forward and the second one backward. The outputs of both LSTMs are connected to one output. This architecture is helpful when the dependencies of a time series go in both time directions, which is the case for snow profiles. (Schuster and Paliwal, 1997; Jurafsky and Martin, 2021)

**Encoder-decoder networks** consist of an ANN encoder that compresses the time-dependent information into a vector and  
155 a decoder that uses this information to solve a supervised sequence labeling task. Additionally, the attention mechanism can be used to strengthen the ability to learn long-term dependencies by focusing only on the parts of the input sequence that are relevant for the current time step. (Bahdanau et al., 2014; Jurafsky and Martin, 2021)

## 2.4 Evaluation

In this work, (1) the performance of different models is compared, (2) differences in the classification of different snow types  
160 are analyzed, and (3) the generalization capability of the best-performing model is examined. (1) The performance comparison is done by looking at the metrics of each model and the specific predictions on the test data set. The metrics used here are accuracy, balanced accuracy, weighted precision, AUROC, log loss, fitting, and scoring time. (2) The label-wise performance is analyzed with the help of label-wise accuracy plots and ROC curves. ROC is the receiver operating characteristic and plots the true positive rate versus the false positive rate. The higher the area under the ROC curve (ROC AUC / AUROC), the clearer  
165 can the model separate between positive and negative samples. (3) The generalization capability is tested by running the best-performing model on 100 random profiles from different parts of MOSAiC winter data. This data is outside of the distribution of the training, validation, and testing data, however, it contains the same classes as the training data, i.e. the model still has a chance to predict the correct labels. Evaluating these three aspects ensures that practitioners can choose a model and know  
170 (1) how it performs compared to other models, (2) what to expect from the snow type specific predictions, and (3) how robust their model will be.



Category	Model	Absolute Accu- racy	Balanced Accu- racy	Prec- ision	F1 Score	ROC AUC	Log Loss	Fitting Time	Scoring Time
Baseline	Majority Vote	0.39	0.14	0.15	0.22	nan	nan	< 1	< 10 <sup>-3</sup>
Semi- Supervised	K-means	0.62	0.44	0.60	0.61	nan	nan	385	0.01
	GMM	0.65	0.36	0.57	0.61	nan	nan	151	<u>0.008</u>
	BGMM	0.65	0.38	0.63	0.63	nan	nan	225	0.009
	Self trainer	0.69	<u>0.67</u>	<u>0.74</u>	0.71	0.92	0.84	19	0.29
	Label propagation	<u>0.71</u>	0.54	0.72	<u>0.71</u>	0.92	1.5	<u>10</u>	3.35
Supervised	Random Forest	<u>0.73</u>	0.60	0.73	<u>0.73</u>	0.93	0.70	72	0.97
	Balanced RF	0.70	<b>0.67</b>	<u>0.74</u>	0.71	0.92	0.84	9.9	<u>0.58</u>
	SVM	0.71	0.66	0.73	0.71	<u>0.93</u>	<u>0.67</u>	19	7.45
	KNN	0.71	0.54	0.71	0.71	0.89	3.58	<u>&lt; 1</u>	1.84
	Easy Ensemble	0.62	0.59	0.70	0.64	0.88	1.66	46	42.5
ANNs	LSTM	<i>0.75</i>	<u>0.58</u>	<i>0.75</i>	<i>0.75</i>	<b>0.94</b>	<b>0.63</b>	<u>349</u>	<u>2.3</u>
	BLSTM	0.74	0.58	0.74	0.73	0.93	0.79	975	3.4
	Encoder-Decoder	<b>0.78</b>	0.54	<b>0.78</b>	<b>0.77</b>	<i>0.94</i>	<i>0.64</i>	2911	5.8

**Table 1.** Results of different models from the categories baseline, semi-supervised, supervised and ANNs. The best values among all models are **bold**. Second-best values among all models are *italic*. The best values among one category are underlined. ROC AUC and logistic loss (log loss) could not be determined for the baseline and some of the semi-supervised models due to the design of these models.

## 2.5 Experimental setup

The experimental setup includes a training, validation, and testing framework: roughly 80% of the labeled dataset is used for training and validation, while the other 20% is set aside for testing. Validation is realized as a 5-fold cross-validation (Stone, 1974). The hyperparameters were tuned on the validation data and the best found hyperparameters were used during testing.

175 Hyperparameter tuning is performed on the validation data. The best found hyperparameters are used for testing. Moderate hyperparameter tuning was applied and all tuning results can be found in the GitHub repository. Specifications of the machine on which the experiments were run can be found in Appendix A and descriptions of the model setup can be found in Appendix B.





### 3 Results

#### 180 3.1 Classification performance of models

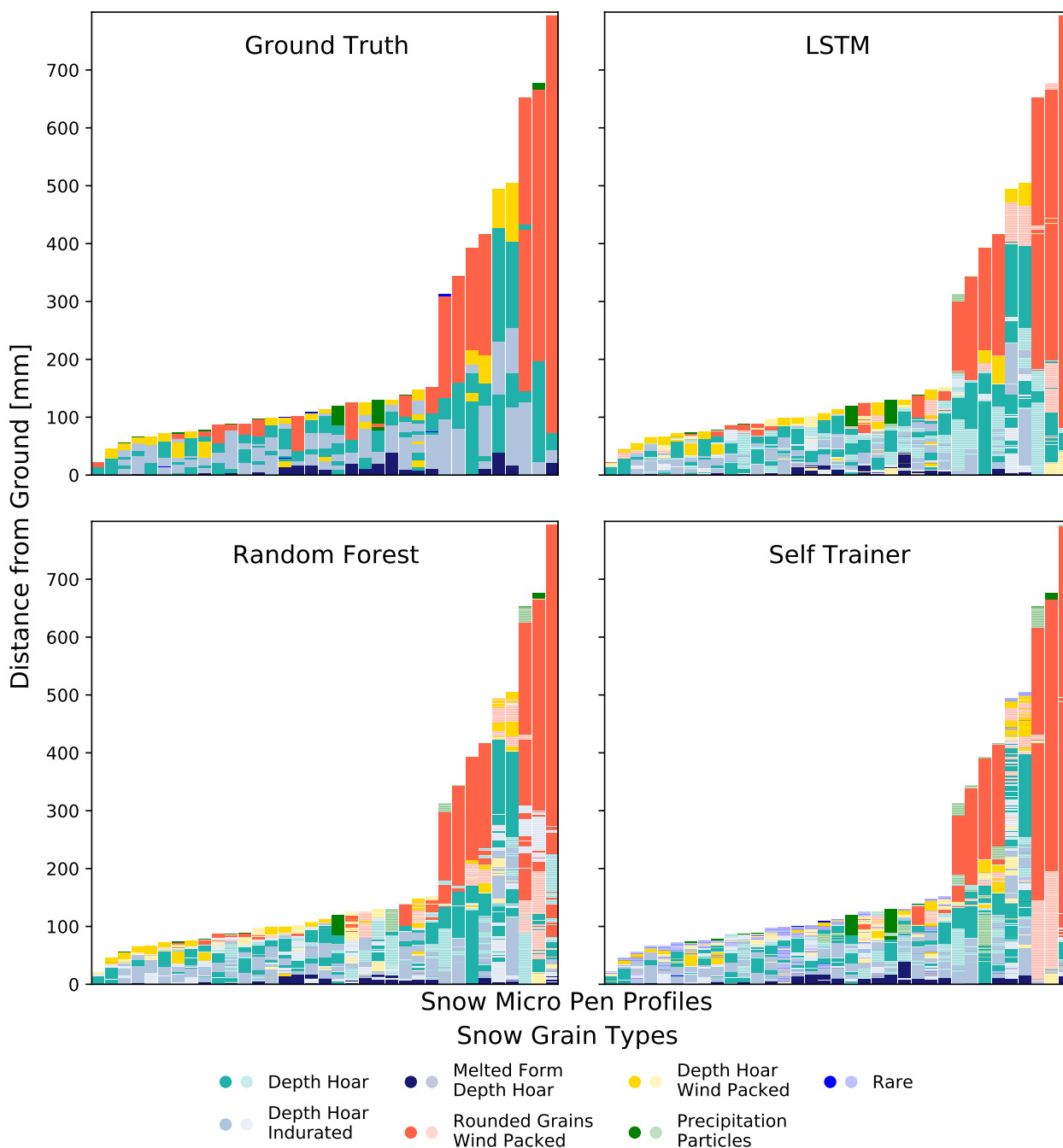
Overall, the results show that an automatic classification and segmentation of SMP profiles with ML algorithms is possible, even if no further information such as snow-pit data or manual segmentation is provided. Category-wise all semi-supervised models were not performing particularly well (see Table 1). Only the self trainer could compete with models from other categories, but this might be the case because the self trainer is based on the balanced random forest. The supervised models achieved mixed performances: Some models such as the random forests and the SVM are clearly performing well, whereas other models such as the KNN and the easy ensemble are underperforming. Overall, the random forest was the best model in the supervised category since it achieves the highest absolute accuracy (0.73) and F1-Score (0.73). However, considering rare classes, the balanced random forest outperformed the plain random forest. All three ANNs did exceptionally well and their category was clearly the most successful among all three categories. The encoder-decoder showed the best scores among all models in terms of absolute accuracy, precision, and F1-Score, closely followed by the LSTM. We consider the LSTM the best model within that category since the encoder-decoder only reached its high performance after extensive hyperparameter tuning and underperformed significantly when not tuned well. In contrast, the LSTM achieved its performance more consistently and even under moderate hyper-parameter tuning, and is thus more suitable for practitioners. The subsequent analyses compare those three models that performed best within their category: the LSTM performed best among the ANNs, the random forest among the supervised models, and the self trainer among the semi-supervised models.

Different ML models exhibited different prediction styles in terms of smoothness and ability to predict rare classes. In Fig. 2 it becomes visible that the models' predictions are not far off from the ground truth. In general, the predictions are somewhat similar to the ground truth but the models often had difficulties in determining the precise start and end of a segment. Looking at three random exemplary profiles of the test data in Fig. 3, one can see that the three main models seem not only to generate similar predictions, but make also similar mistakes. In the medium-deep profile (middle column), all three models predicted a longer segment of "Depth Hoar" that was actually not present in the ground truth profile. In the shallow profile all three models predicted some intermediate "Depth Hoar Wind Packed" layers in the first third that did not exist. And in the deep profile, all three models miss the narrow intermediate "Depth Hoar" layer. In summary, it becomes apparent that the different models are producing consistent predictions to a certain degree. Of course, there are significant differences among the models, too. First of all, the LSTM is closest to the ground truth (see Fig. 3). Secondly, the LSTM provided much smoother and less fragmented predictions than the other two models. And thirdly, the self trainer clearly overestimates rare classes, which hurts the overall performance. To summarize, the LSTM, random forest, and self trainer show certain prediction similarities among each other, however, the LSTM is closest to the ground truth and imitates expert labeling best.

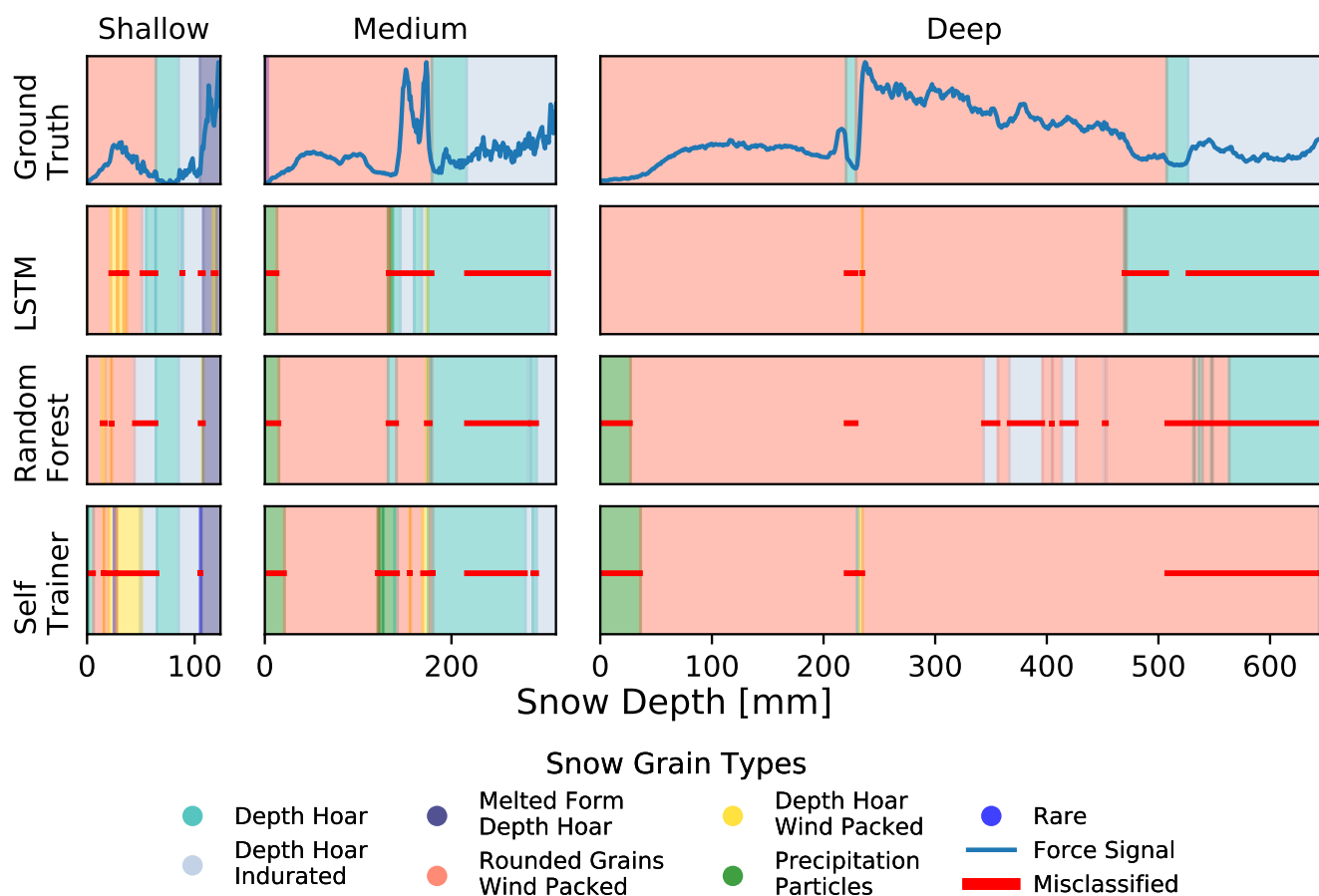
#### 3.2 Classification difficulty of snow types

210 Fig. 4 shows that some snow types are easier and others are harder to classify. The label-wise accuracy seems to be influenced by the following factors: (1) choice of model, (2) frequency of snow type in the dataset, (3) snow type itself. Within one snow





**Figure 2.** Predictions on the test dataset of the LSTM, random forest, and self trainer. The upper left panel shows the ground truth data. In the other panels, the correct predictions are shown with more intense colors and the wrong predictions with less intense colors. The LSTM has the highest rate of correct predictions and imitates the smoothness of the ground truth very well. The random forest does well but provides more segmented predictions. The self trainer immensely overestimates rare classes.



**Figure 3.** Model predictions for three randomly chosen SMP profiles. The first row represents the ground truth labels (with force signal). The subsequent rows represent the LSTM’s, random forest’s, and self trainer’s predictions, with the red bar indicating wrong predictions. Each column shows a different profile randomly chosen from the test data (shallow profile: S31H0276; medium profile: S31H0206; deep profile: S49M1918). All three models seem to make similar mistakes, e.g. they predict a larger portion of “Depth Hoar” at the end of the medium SMP profile. The predictions of the LSTM are closest to the ground truth data.



type category, the models perform differently well, however, some snow types seem to be easier, and other more difficult to classify for all models. For example, “Rounded Grains Wind Packed” achieved a high accuracy among all models, whereas “Depth Hoar Wind Packed” achieved a low accuracy among all models. This could be partially attributed to the fact that there are fewer samples available for “Depth Hoar Wind Packed”. However, the snow types themselves seem to influence the classification difficulty as well: the class “Precipitation Particles” achieves high accuracy values among some models, despite the fact that it is the rarest class in the dataset. For some snow types, some models are able to access certain information enabling a high performance on that particular snow type – independent of its frequency. This means that the classification difficulty does not only depend on the number of available samples, instead, some other underlying characteristics determine the classification difficulty of the snow types as well.

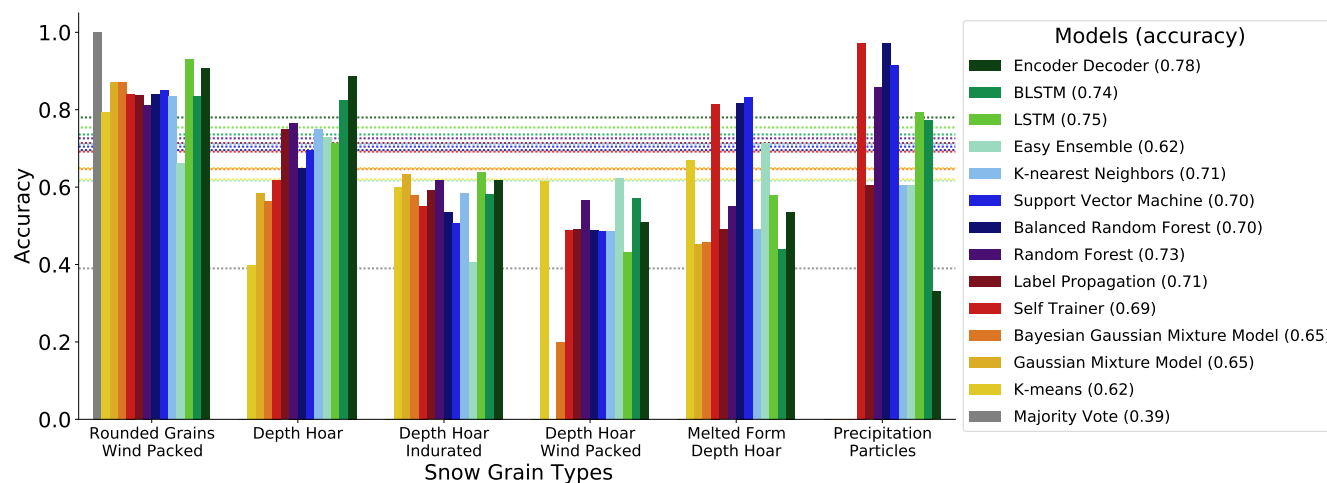
Depending on the model, a higher accuracy score could lead to a lower precision score for a label (accuracy-precision trade-off). The ROC curve in Fig.5 illustrates this relation between the true positive and false positive rate for the different snow types and their averaged performances. It becomes apparent that both the snow type and the choice of model influence the accuracy-precision trade-off. The class “Rare” for example seems to be difficult to classify both accurately and precisely for all models, whereas “Precipitation Particles” are showing an almost perfect ROC curve. If one is interested in choosing a model that performs well for a particular snow type, these ROC curves can reveal which model is most suitable. To get even more detailed label- and model-wise insights, refer to the confusion matrices in Appendix E. Both the LSTM and the random forest achieve an area under the ROC curve of 0.96. However, on average (see Fig. 5, pink dotted line), the LSTM outperforms the self trainer and random forest and is thus most suitable for general classification tasks.

### 3.3 Generalizability

The prediction of the LSTM for 100 random profiles outside of the training and testing distribution is shown in Fig. 6. Since the ground truth profiles are not yet available for these predictions, the generalization capabilities can only be evaluated on the basis of what seems “reasonable”. “Melted Form Depth Hoar” appears only at the ground of the profiles, “Precipitation Particles” only at the top, “Rounded Grains Wind Packed” are mostly at the top and rather deep – these are all “reasonable” predictions. However, there are also some predictions that are not reasonable or at least unexpected: the left profile consists almost entirely of “Depth Hoar Wind Packed”, sometimes “Depth Hoar Wind Packed” appears right before “Melted Form of Depth Hoar”, and “Rounded Grains Wind Packed” sometimes appear briefly in the “middle” of a profile (and not at the top). Overall, the LSTM seems to make mostly reasonable predictions, however, an in-depth expert analysis of the predictions is necessary to validate that further.

## 4 Discussion

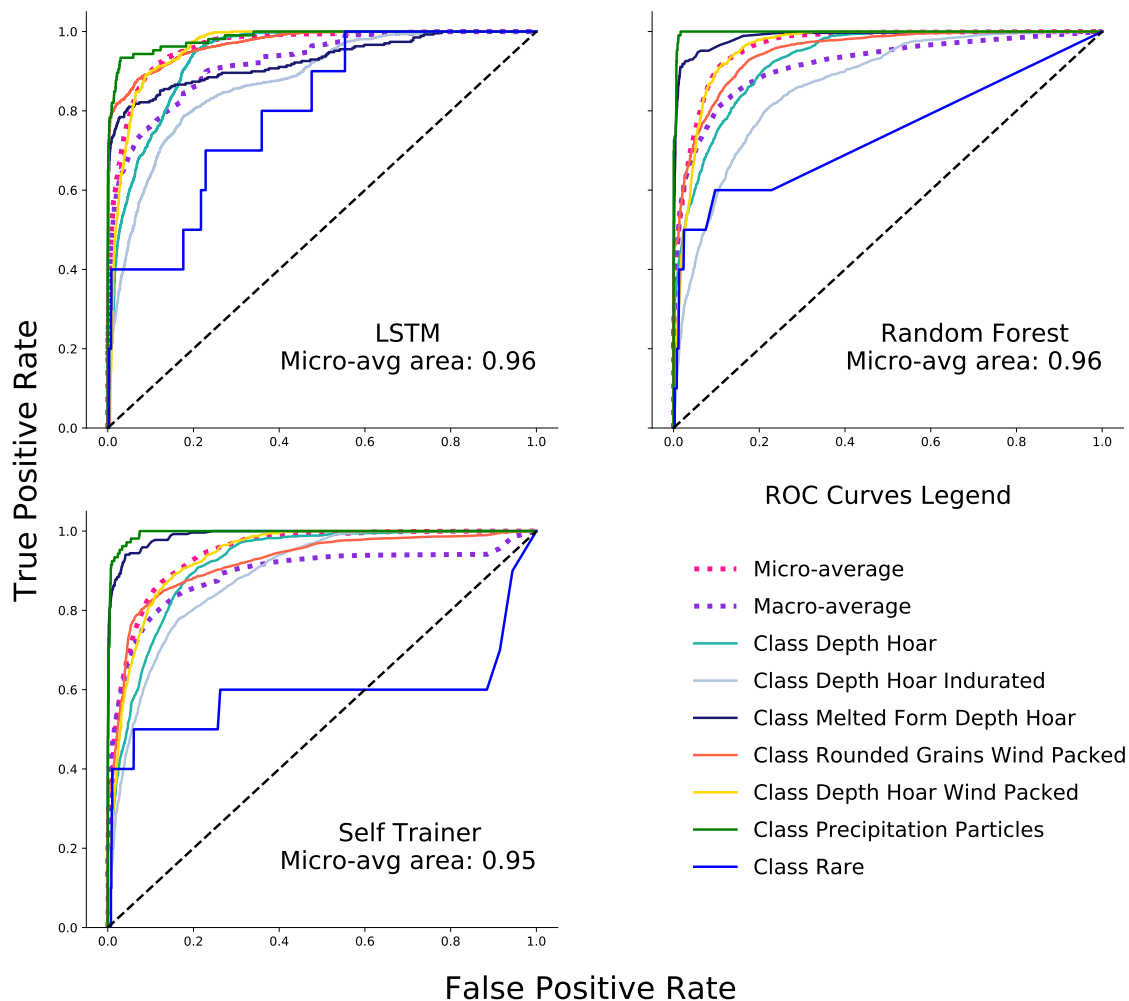
The results showed that automatic classification and segmentation of SMP profiles is possible with up to 78% accuracy. In the following the nature, impact, and limits of these results are discussed.



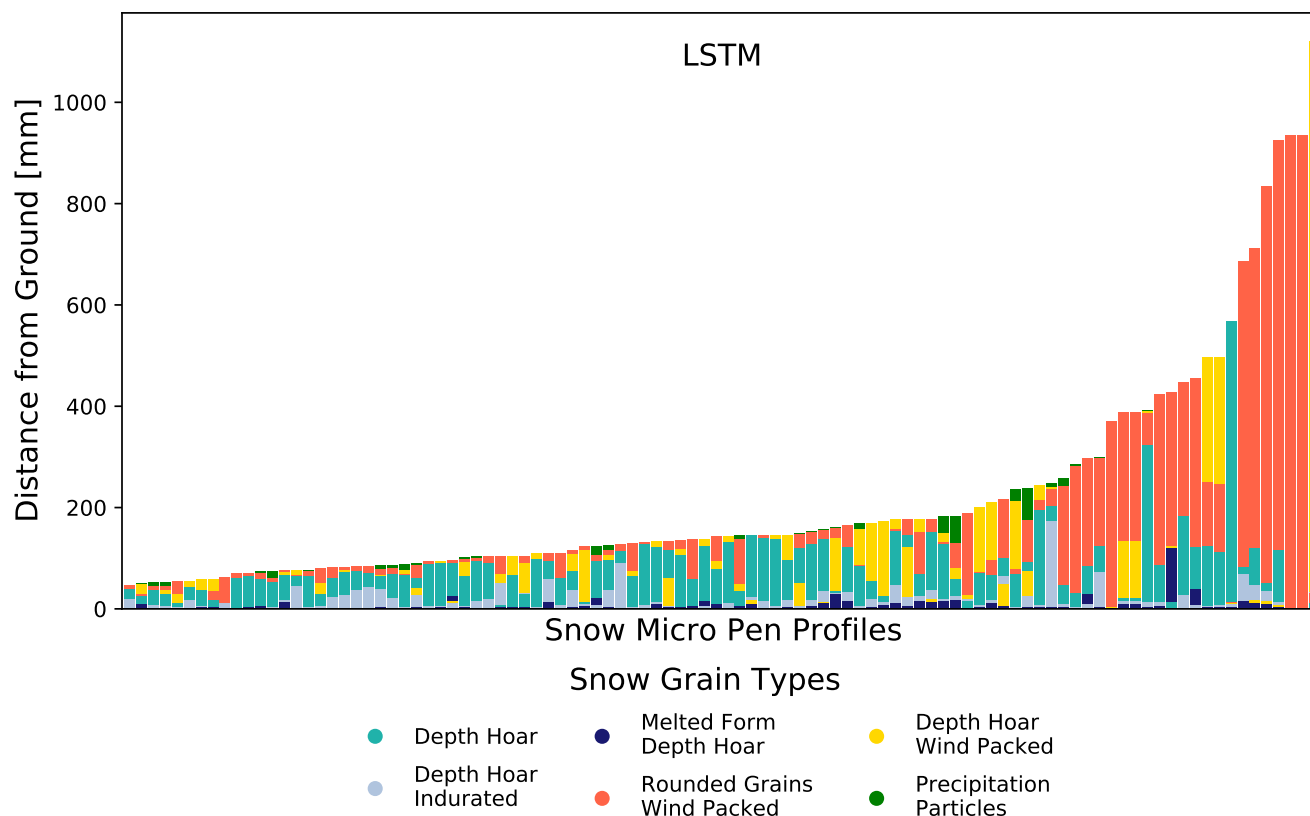
**Figure 4.** Label-wise accuracy of all models. each model is encoded with a different color. The most frequent label is on the left of the x-axis (“Rounded Grains Wind Packed”), and the least frequent on the right (“Precipitation Particles”). The class “Rare” was dropped. Each bar represents the accuracy for a single snow type. The dotted lines show the overall accuracy performance of each model. The encoder-decoder, the BLSTM, and the LSTM achieved the highest accuracy values. For all models, some classes are more difficult to classify than others: e.g. “Depth Hoar Indurated” and “Depth Hoar Wind Packed”. Some classes are easier to classify than others, such as “Rounded Grains Wind Packed”. Some classes can only be classified well by a subset of the models, such as “Precipitation Particles” and “Melted Form Depth Hoar”.

The metrical results presented are in line with previous findings: King et al. (2020) reported an overall accuracy score of 0.76 when using SVMs and additional snowpit information to classify three snow types. Satyawali et al. (2009) achieved an average accuracy of 0.81 when using the nearest neighbor approach and knowledge rules to classify five snow types. However, these results stem from only three profiles and are not representative. Havens et al. (2012) achieved an accuracy of maximal 0.76 (global dataset) when using random forests and time-intensive manual layer segmentation to classify three snow types. The major difference from these previous results is that the accuracy results of this study were achieved for *seven* snow types, without time-intensive layer picking, snowpit digging, or additional knowledge rules. This means that in contrast to previous work, the models here can be directly employed by practitioners for their own SMP datasets in the field: simply retrain and predict. For this, they only need to provide a set of training samples for their specific dataset and classification style. The work presented here enables scientists for the first time to rely on fully automated ML SMP profile segmentation and classification.

The results were also satisfying to domain experts since the predictions were in themselves consistent and followed the patterns of the training data. In general, the snowpack on sea ice is extremely variable, and the traditional snow types are very often a mixture of different features. This becomes visible when comparing the SMP-profiles to the micro-CT samples. In the view of the authors, a temporally consistent classification is more relevant to the interpretation of the development of the snowpack, even if there is a certain, but unknown, bias to an expert interpretation. Hence, the models were also in practice helpful to analyse Arctic snowpack development.



**Figure 5.** ROC curves of the LSTM, random forest, and self trainer for each class. The dotted lines are the micro- and macro-averaged ROC curves. The macro-average calculates the ROC for each class and averages the performances afterward. The micro-average weights the performance according to class contribution (balanced performance results). The LSTM achieves the highest ROC performance overall. The order of the best-performing snow types is similar among all models. The classes “Rare” and “Depth Hoar Indurated” have the lowest ROC areas, whereas “Precipitation Particles” has the highest ROC area for all models.



**Figure 6.** LSTM SMP profile predictions on out-of-distribution data. The SMP profiles used here come from different legs of the MOSAiC expedition than the training, validation, and test data. The profiles used here still stem from the winter season to ensure that the same set of snow types can be used as in the training dataset. The distribution of the predicted profiles looks convincing, with only a few profiles standing out as certainly wrong predictions (e.g. most right profile with  $\sim 90\%$  “Depth Hoar Wind Packed”).

#### 4.1 Classification performance of models

260 Each model category are performs differently because each model takes different aspects of the data into account. Semi-supervised models try to take unlabeled data into account to improve their predictions, however, this did not work well in our context. The most likely reason for the overall underperformance of this category is that the unlabeled data contained out-of-distribution data, i.e. the unlabeled data had different underlying mechanisms than the labeled data (different parts of the winter season). Another reason might be that only a small subset of unlabeled data was included in order to limit running  
265 times. Moreover, the poor performance of the cluster-then-predict-models is most likely also a result of the classifier used after clustering: a more sophisticated method than a majority vote classifier is needed here.

The simple supervised models take one data point after the other into account and do not consider time-series structures within the data. The algorithms used in all previous SMP automation studies fall into this category. In contrast, ANNs are



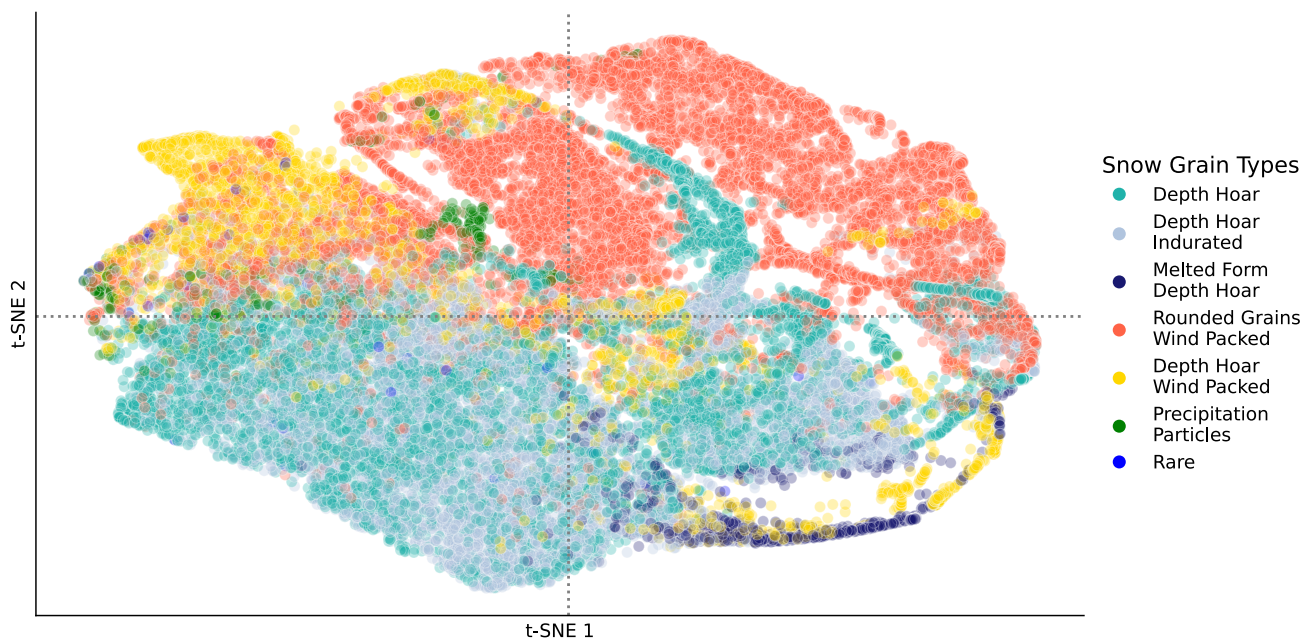
270 supervised models that take the underlying time sequence of the data into account. While the supervised model in general  
performed well, they were still clearly outperformed by the ANNs. A likely reason why the ANNs outperformed all the other  
models is precisely the ANNs' ability to process time-dependent information. ANNs are tackling the classification task as a  
sequence labeling task which enables them to include information from the order and position of snow layers. The supervised  
models still have access to time-relevant information (time-window features), however, they do not have any ability to learn  
time-based information (what should be remembered and forgotten). Besides, the ANNs learn to imitate the training set, leading  
275 to smooth and expert-simile predictions. In comparison, taking the time component of SMP signals into account has not been  
done in previous methods and we argue that it adds a major information piece and boosts the overall prediction performance  
significantly.

Each model exhibits a different prediction style due to the models' intrinsic differences and thus might be suitable for specific  
tasks. In the following some aspects are listed for consideration (practitioner's guide):

- 280 **A Time and resources for hyperparameter tuning.** The LSTM and the encoder-decoder network are recommended when  
plenty of tuning time is available. Especially, the encoder-decoder network performs badly if not tuned well. The SVM  
and the balanced random forest need little tuning time, whereas the random forest is the go-to-model in case (almost) no  
tuning time can be provided.
- B Need for a simple to handle, off-the-shelf algorithm.** Among the high-performing models, the random forest and the  
285 SVM are the easiest to handle off-the-shelf algorithms. The self-supervised algorithms and especially the ANNs require  
a somewhat deeper understanding of the models and the ability to implement them.
- C Desired level of explainability.** The random forests are most explainable since the decision trees can be directly visual-  
ized (Appendix D). The ANNs are the least explainable models (without further modifications).
- D Importance of minority classes.** When deciding on a model, the underlying task must be examined as well: In the case  
290 of avalanche prediction it might be essential to predict a buried layer of "Surface Hoar", a very rare class, which needs  
to be detected no matter the costs. In such a case of "minority class prediction" the balanced RF or the SVM should be  
employed. The ANNs and the random forest, in contrast, are more suitable to achieve an overall good classification.
- E Availability of unlabeled data that is from the same distribution as the labeled data.** In case a lot of unlabeled  
data from the same distribution and time is available, the self-trained classifier can be considered. The weak learner of  
295 the self-trained classifier can be chosen according to the criteria listed above. Since in this work we only had a small  
subset of unlabeled data stemming from the same distribution as the labeled data, further evaluations on the self-trained  
classifier and label propagation remain open.

This highlights that there is not a single best model, but instead, practitioners can deliberately choose a model that suits their  
needs, such as overall accuracy, ability to predict rare classes, explainability, training, and deployment time.





**Figure 7.** 2-dimensional t-distributed stochastic neighbor embedding (t-SNE) of SnowMicroPen (SMP) dataset. The colors encode the snow types. The figure shows that (1) “Depth Hoar” and “Depth Hoar Indurated” are hardly separable, (2) “Depth Hoar Wind Packed” is similar to several other snow types, and (3) “Precipitation Particles”, “Melted Form of Depth Hoar” and “Rounded Grains Wind Packed” can each be separated more clearly from the other snow types.

#### 300 4.2 Classification difficulty of snow types

Snow types are differently difficult to classify since their categories are rather continuous than discrete. This was also observed in previous work and in all previous works performances were reported label-wise to account for those differences (Satyawali et al., 2009; Havens et al., 2012; King et al., 2020). We performed t-distributed stochastic neighbor embedding (t-SNE) on the SMP dataset to visualize how separable the different classes are (see Fig. 7). “Precipitation Particles”, for example, appears as a singled-out green island, which is in line with our and other findings (Satyawali et al., 2009) that it is easier to classify than other snow types. We conclude from this, that some classes have features that distinguish themselves stronger from other snow types. The class “Rounded Grain Wind Packed” behaves similarly (Satyawali et al., 2009). However, some classes, such as “Depth Hoar” and “Depth Hoar Indurated” are completely overlapping in Fig. 7, and indeed our models had problems with differentiating between those two classes. Similarly, “Depth Hoar Wind Packed” seems to overlap largely with “Rounded Grains Wind Packed” and “Melted Form of Depth Hoar”. We theorize that the reason for their non-separability is that those snow types transform into each other during snow metamorphosis. This means many data points can not be discretized into one single category since they are on a continuous spectrum. Satyawali et al. (2009) pointed out, as well, that they often found data points being in transition between snow classes and attributed it to the fact that the snow is changing continuously. In



conclusion, it is virtually impossible to reach 100% classification accuracy on every snow type since some snow types will  
315 always lie between two categories.

The classification difficulty of the different snow types extends also to the expert labeling process itself. The continuous  
natures of the labels and additional challenges such as between-class imbalances, make it particularly difficult for domain  
experts to label the SMP profiles consistently among each other. The uncertainty during labeling is an intrinsic problem of  
SMP analysis and cannot be circumvented: The annotation of SMP profiles stays always subjective, meaning that two different  
320 snow experts may produce two different labeled and segmented profiles for the exact same measurements (Herla et al., 2021).  
However, both experts might agree that both labeled profiles are valid analyses of the same profile. Hence, the model's perfor-  
mances cannot only be measured in terms of accuracy because models with low accuracy might still produce sensible, directly  
usable predictions. Throughout our experiments, some of the models' predictions –to our surprise –already helped domain  
experts to detect mistakes and inconsistencies in their ground truth labeling. Due to the experts' individual classification styles,  
325 the models must adapt to those styles to truly satisfy the needs of practitioners. This means the models must be re-trained on  
a data set of the particular practitioner. Alternatively, the models could be used to support and speed up the manual labeling  
process by making label suggestions that are then checked by a snow expert.

### 4.3 Generalizability

The LSTM can generalize to other winter profiles with the same snow types since the underlying classification and segmen-  
330 tation rules stay the same. However, the LSTM's generalization capability does not extend to other seasons or regions when  
/ where other snow types are found, such as melted forms or regional snow types. As mentioned before, the models do not  
generalize on different classification styles of experts. The models used in this work are still generalizable in that they can be  
used on any desired dataset as long as they are re-trained on the chosen dataset. This would not have been possible in previous  
works such as Satyawali et al. (2009) since knowledge rules for one snow region and season do not transfer to other regions or  
335 seasons. For greater generalization capability the LSTM – or any other model — must be either trained with a more general  
dataset or must be specifically re-trained for an individual data set.

### 4.4 Limitations and Future Work

This work does not address the task setting of first-segment-then-classify because this would require a completely different  
set of methods. In a first-segment-then-classify setting, the SMP signal could first be segmented with techniques used in  
340 audio-segmentation (Theodorou et al., 2014). The resulting time-series pieces could subsequently be classified as a whole  
(Ismail Fawaz et al., 2019). Future work could experiment with this problem formulation and analyze if performance further  
increases in this setting.

The ANNs used here are off-the-shelve and are not adapted to the specific underlying task in order to ensure a fair comparison  
between the different models. However, one could look into adapting the loss functions to include similarity measurements  
345 between snow samples. Results from clustering, performed on t-SNE data, could then be leveraged during classification to



increase classification performance. Adapting the loss function of the ANNs could increase prediction performance greatly, however, such a loss function must be carefully constructed and evaluated on different datasets.

As mentioned in Sect. 4.3, the models cannot generalize to completely different settings in terms of seasons and regions. To ensure generalization capability one could train a large model on a dataset that includes snow types from different regions and seasons. Such a data set would need to be newly compiled because common SMP datasets are usually limited to one region (Ménard et al., 2019; Calonne et al., 2020). However, it is completely unclear if classification on such a large dataset would actually yield better performances. The classification task does become significantly harder with more classes and different data distributions. However, a large enough model trained on a large enough dataset could in theory be able to produce direct predictions for any SMP user. Alternatively, SMP users can simply re-train a chosen model for their particular dataset. They would need to provide a set of SMP profiles for their region, season, and classification style, but the overall time savings are still immense. To summarize, the generalization capabilities may be enhanced by using a more general dataset or one bypasses this problem by re-training to specific datasets – the snowdragon repository addresses the needs of the latter.

An immediate consequence of this study is the further analysis of the unlabeled part of the MOSAiC dataset. Domain experts can use the LSTM, or other models, to create predictions for the remaining 3516 profiles. A previously almost impossible task to classify and segment those thousands of profiles, became feasible by providing just a set of 164 labeled profiles. The results of these predictions and their impacts on the cryospheric analysis of snow coverage in the Arctic will become apparent in future publications.

## 5 Conclusions

This study showed for the first time that SMP profiles can be automatically segmented and classified (up to 0.78 accuracy). Fourteen different models were trained here to classify seven snow types without providing any additional manual information. It also showed for the first time how ANNs and semi-supervised models can be used for the task of SMP classification and segmentation. Among all models, the LSTM and the encoder-decoder are performing the best. The resulting predicted profiles show smooth segmentations and expert-simile classification patterns that were satisfying to domain experts.

These findings will enable SMP practitioners to automatically analyze their SMP measurements. To that end, an SMP user must simply decide on one of the fourteen models provided by the snowdragon repository, given the considerations listed in this paper, and retrain the model for their particular dataset. Afterward, the SMP user can simply predict SMP classifications and segmentations for the remaining unlabeled profiles.

Snowdragon could be extended further, made more user-friendly, and in particular, it could be integrated into the snowmicropyn package. The resulting tool would make knowledge about snowpacks easier and faster accessible for all scientists. This is of particular interest (1) for interdisciplinary scientists who rely on snow type but do not have the tools to classify them themselves (remote sensing), (2) for scientists that require fast analysis of SMP profiles, such as in avalanche prediction and (3) for SMP users facing large datasets.



Snowdragon enables already today the analysis of the SMP MOSAiC dataset with a large amount of detailed data about the Arctic's condition. The ML-driven approach used here to analyze SMP profiles will be one of many methods to make the knowledge behind the data accessible – knowledge that is essential to understanding and mitigating climate change impacts.

*Code and data availability.* The current version of snowdragon is available on GitHub: <https://github.com/liellnima/snowdragon> under the MIT licence. To run the code version used in this paper, please refer v1.0.0 on GitHub or Zenodo: <https://doi.org/10.5281/zenodo.7335813>. The exact version of the models used to produce the results used in this paper is also archived on Zenodo: <https://doi.org/10.5281/zenodo.7063520> (Kaltenborn et al., 2022). The MOSAiC SMP data used as input and training data is available on PANGAEA: <https://doi.pangaea.de/10.1594/PANGAEA.935554> (Macfarlane et al., 2021).

## Appendix A: Machine specifications

The evaluation and hyperparameter tuning experiments were run on two different machines. The complete evaluation was conducted on a 64-bit system with an Ubuntu 18.04.5 (Bionic Beaver) operating system. The machine has 16 GB RAM and an Intel® Core™ i7-6700HQ CPU @ 2.60GHz × 8 (and the GPU was not used). The machine on which the first hyperparameter tuning, training, and validation experiments have been run has the following specifications: 64-bit system with an Ubuntu 20.04.1 (Focal Fossal) operating system, an Intel® Core™ i7-4510U CPU @ 2.00GHz x 4 CPU, and 12 GB RAM (and the GPU was not used). Final hyperparameter tuning, training, and validation (results presented here) were run on an Azure virtual machine of the Dsv3-series, namely on a Standard\_D4s\_v3<sup>2</sup> machine with Ubuntu 18.04 (Bionic Beaver) as an operating system, 16 GB RAM and 4 vCPUs.

## Appendix B: Model setup

The project was executed in Python 3.6 and all used packages can be found on GitHub in the “requirements.txt” file. Principle component analysis, t-SNE, k-means clustering, Gaussian Mixture Models, Bayesian Gaussian Mixture Models, random forests, SVMs, and the k-nearest neighbor algorithm were used as made available through scikit-learn by Pedregosa et al. (2011).<sup>3</sup> The easy ensemble for imbalanced datasets and a balanced variant of the random forest are imported from imbalanced-learn by Lemaître et al. (2017).<sup>4</sup> All ANN architectures were created with the help of TensorFlow (Abadi et al., 2015)<sup>5</sup> and Keras (Chollet et al., 2015)<sup>6</sup>. The attention model within the encoder-decoder network was used as provided in the keras-attention-mechanism package by CyberZHG (2020).

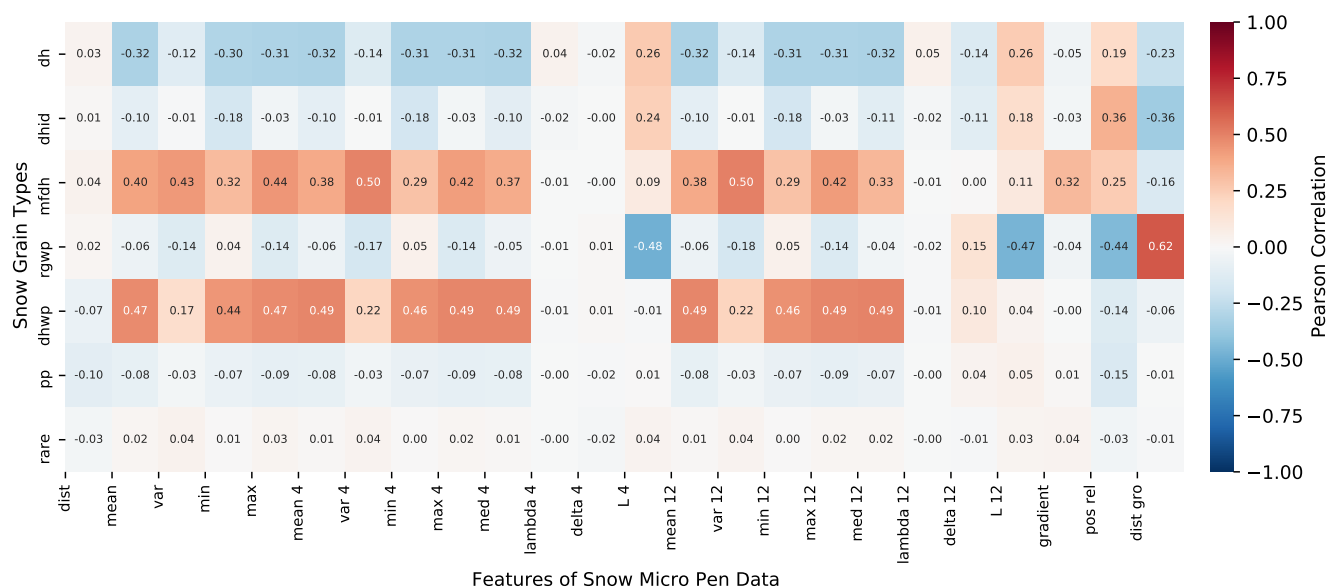
<sup>2</sup><https://docs.microsoft.com/en-us/azure/virtual-machines/dv3-dsv3-series>

<sup>3</sup><https://scikit-learn.org/stable/>

<sup>4</sup><https://imbalanced-learn.org/stable/>

<sup>5</sup><https://www.tensorflow.org/>

<sup>6</sup><https://keras.io/>

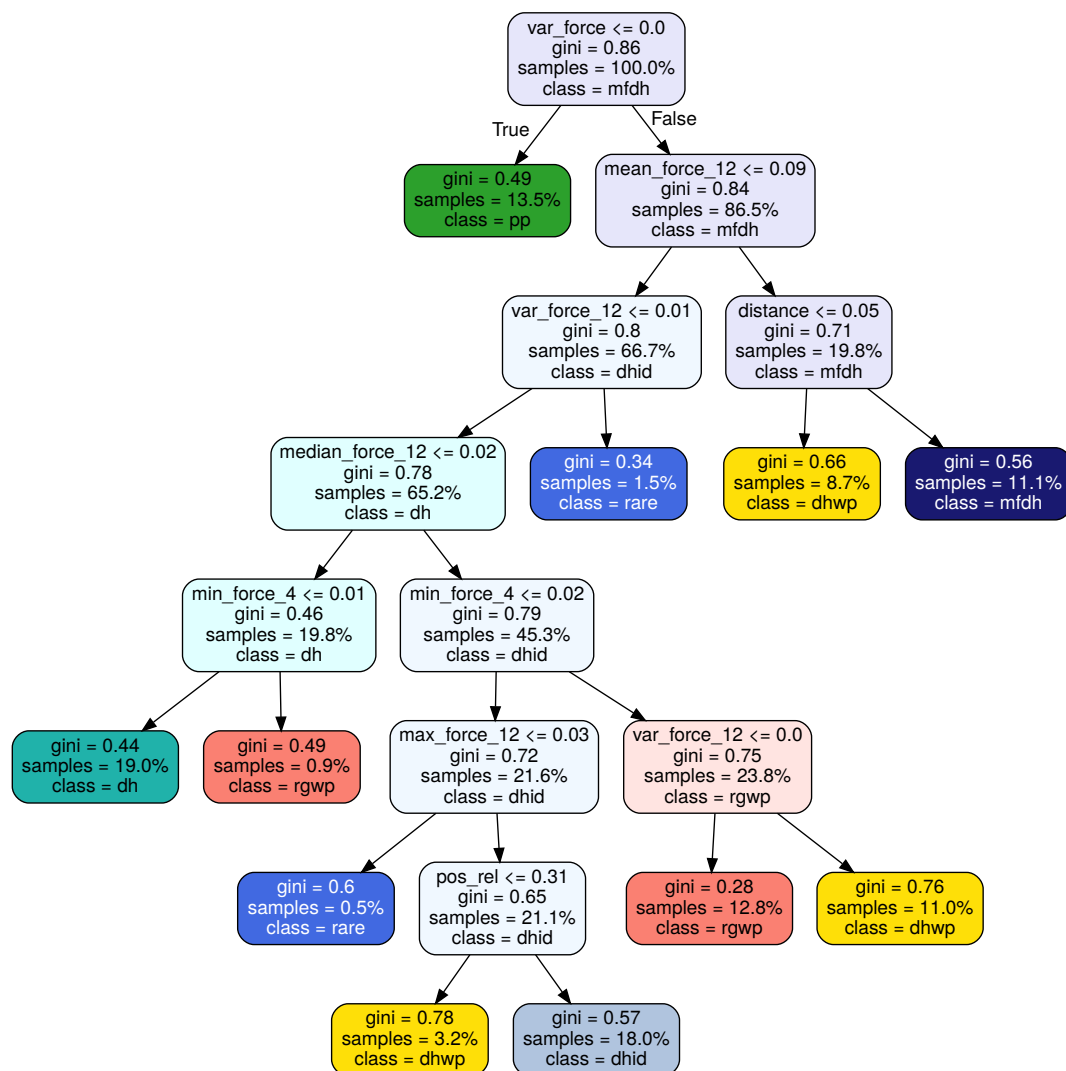


**Table C1.** Label-Feature correlation between snow types and aggregated features of the SMP profiles. The numbers in the feature names stand for the window size used during aggregation. “Depth Hoar” (dh), “Depth Hoar Indurated” (dhid), and “Rounded Grains Wind Packed” (rgwp) show some negative correlations with a subset of the features. “Melted Form of Depth Hoar” (mfdh), “Depth Hoar Wind Packed” (dhwp) and “Rounded Grains Wind Packed” (rgwp) show a strong positive correlation with at least one feature. “Precipitation Particles” (pp) does not show strong correlations with any feature, however, a correlation with distance (dist), variance, and force features was expected by experts. The low correlations could be caused by the data-preprocessing step when “Decomposed and Fragmented Precipitation Particles” were categorized as “Precipitation Particles” as well. The class “Rare” shows no correlations with the features since it consists of very different sub-classes (“Ice Formation” and “Surface Hoar”).

### Appendix C: Label-wise feature correlation

Table C1 shows why classification for this dataset is so hard. Some labels have lower correlations among all features, making it unclear how the right predictions can be achieved on this basis. Other more predictive features are missing, i.e. if a feature is discovered that shows a high correlation within this plot, it might boost the overall classification capabilities of the models. The figure also shows that there might be interaction effects arising since some snow types show very similar correlations (for example “Melted Form of Depth Hoar” and “Depth Hoar Wind Packed”). In summary, the label-wise feature correlation reveals the classification difficulty of the dataset and can be used to discover new predictive features.

### 410 Appendix D: Pruned decision tree

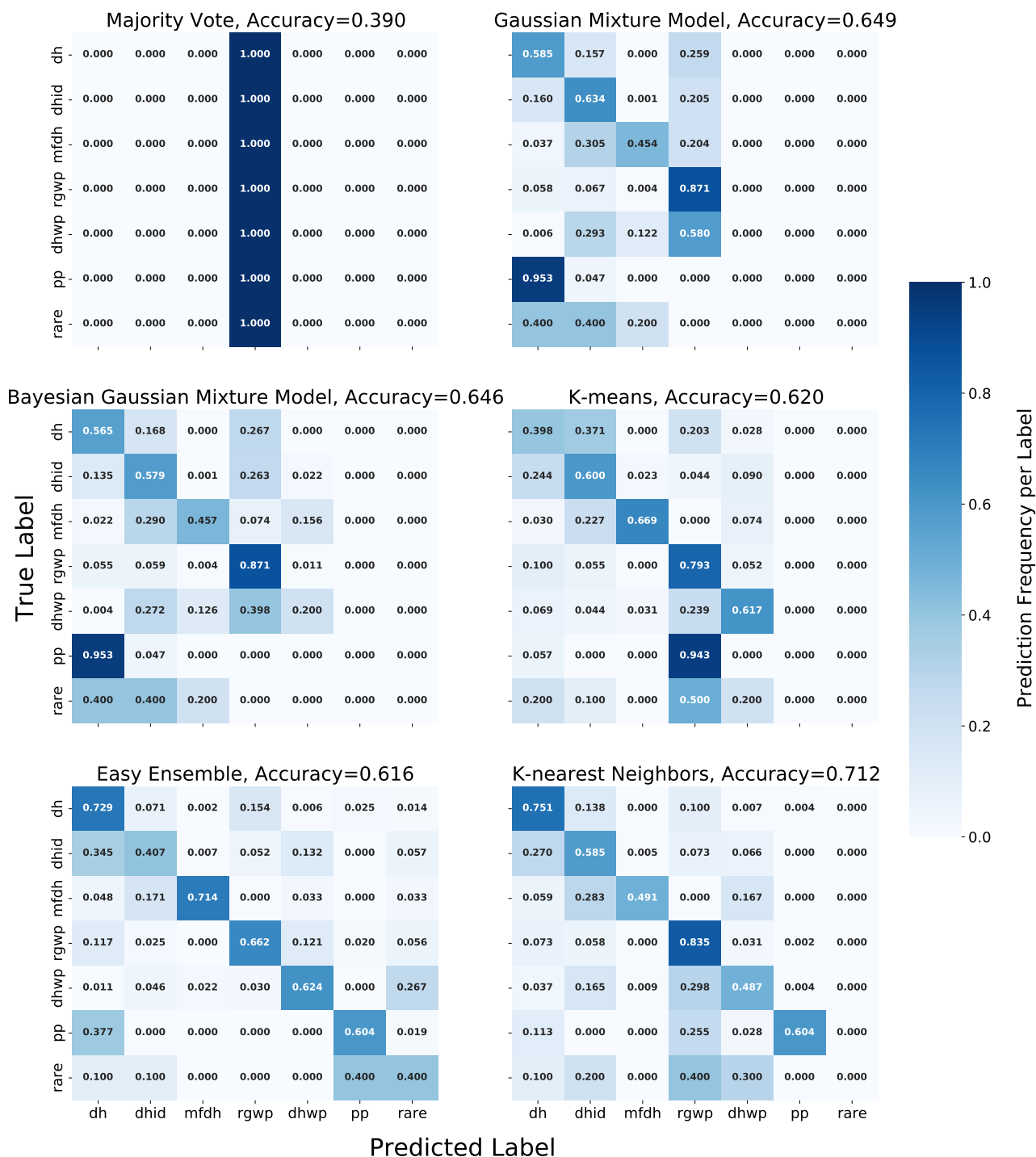


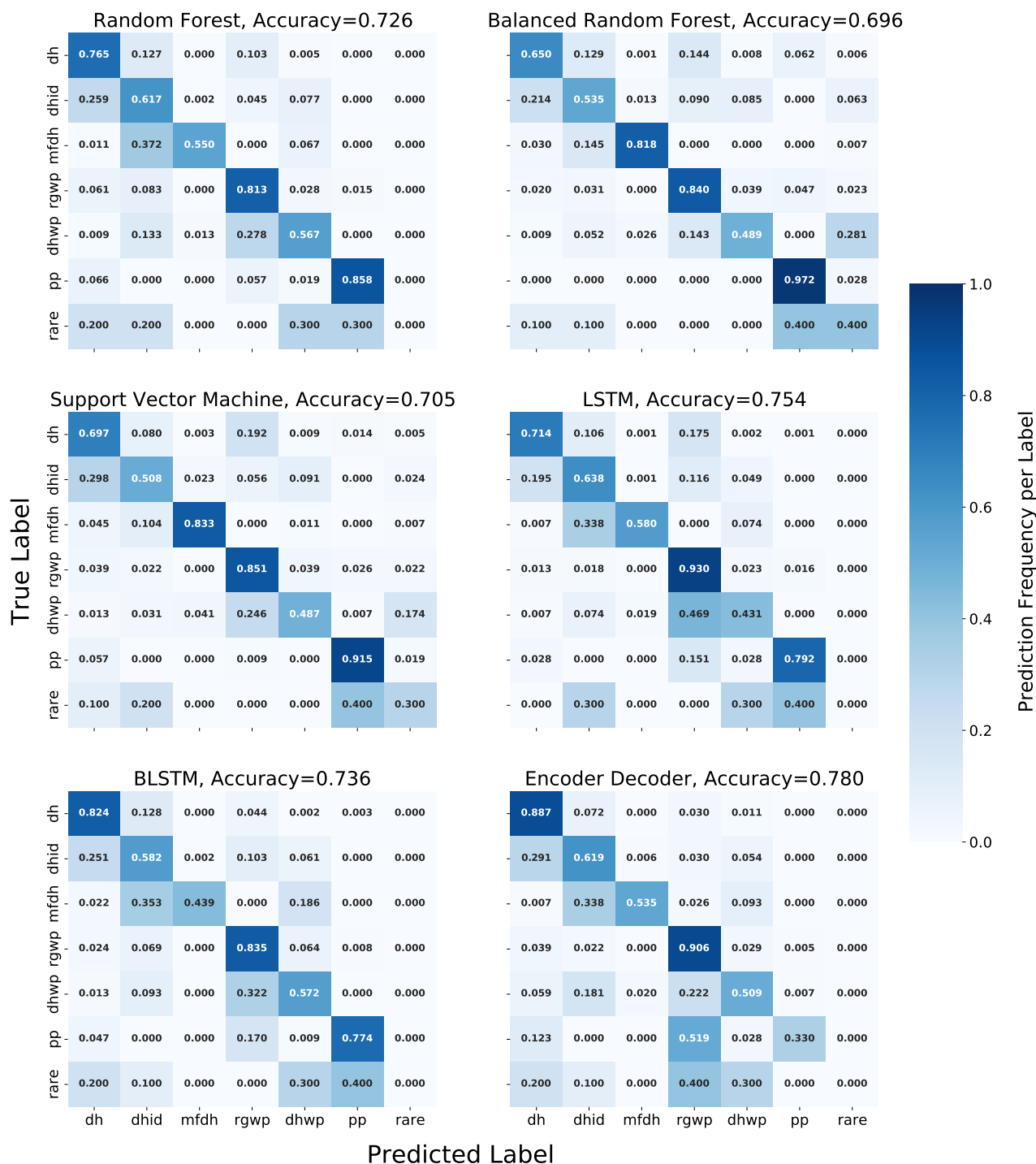
**Figure D1.** Pruned decision tree extracted from the random forest. Decision trees encode the decision rules for predicting snow type labels. This approach helps to explain the model’s decisions, a property that is often asked for by domain experts. At each leaf node, a labeling decision is made. All the other nodes encode the labeling rules that are used to classify each point. Take the root node as an example: If the variance of the force is smaller or equal to zero, the point is labeled as “Precipitation Particles”. Else it has to be one of the other labels. The Gini index encodes how well separable the subsets of data points are (the bigger the number the better), and the sample’s number shows how much percent of the complete data can be found in this subset.

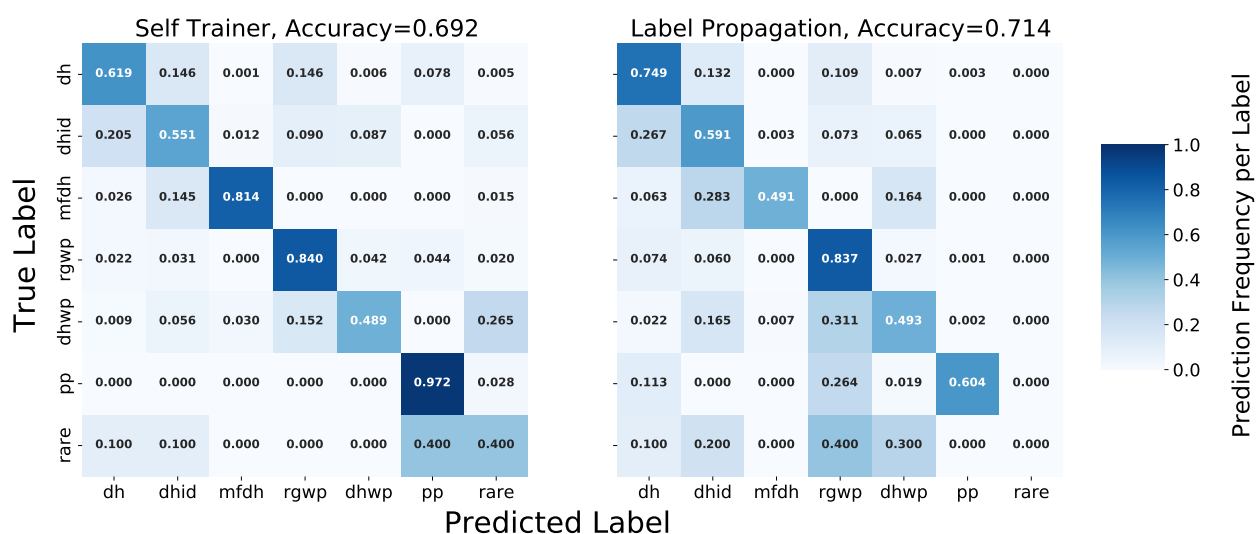


## **Appendix E: Confusion matrices**









**Table E1.** Confusion matrices of all models displaying the predicted and the observed snow types. The number in each cell is the relative prediction frequency of a label within the observed class. The numbers of the diagonal (upper left to lower right) represent the prediction accuracy of each label. The stronger pronounced the diagonal and the less pronounced the upper and the lower triangles are, the better are the predictions. The confusion matrices help for an in-depth analysis of the label-specific performances. This is useful when practitioners want to choose a model that is suitable for a specific snow classification task.



*Author contributions.* ARM and MS collected and curated the data; ARM labeled the data; ARM and JK preprocessed the data; JK developed the methodological framework; JK implemented, compared, tuned and validated the models; JK and VC visualized the results; JK wrote the manuscript draft; VC, ARM and MS reviewed and edited the manuscript; VC supervised the ML part of the study; MS supervised the cryospheric part of the study.

*Competing interests.* The authors declare that they have no conflict of interest.

*Additional notes.* The manuscript might have some similarity with a paper that we submitted to the Climate Change AI Workshop at NeurIPS 2021 (<https://s3.us-east-1.amazonaws.com/climate-change-ai/papers/neurips2021/48/paper.pdf>). The paper submitted to Climate Change AI was a preliminary version of the submitted manuscript and was not peer-reviewed (only superficially checked for scientific correctness). Specifically, snow specific information is only summarized there. The workshop organization committee states on their website (<https://www.climatechange.ai/events/neurips2021>): "The workshop does not publish proceedings, and submissions are non-archival. Submission to this workshop does not preclude future publication."

*Acknowledgements.* This project was funded by the Swiss Polar Institute (DIRCR-2018-003), the European Union's Horizon 2020 research and innovation program projects ARICE (grant 730965) for berth fees associated with the participation of the DEARice project, the WSL Institute for Snow and Avalanche Research SLF (WSL\_201812N1678). The project was additionally financed by the funds of a research training group provided by the Deutsche Forschungsgemeinschaft (DFG), Germany (GRK2340). Data used in this manuscript was produced as part of the international Multidisciplinary drifting Observatory for the Study of the Arctic Climate (MOSAiC) with the tag MOSAiC20192020. The data was collected during the Polarstern expedition AWI\_PS122\_00. We acknowledge the contribution of the MOSAiC-expedition (Nixdorf et al., 2021). We especially thank the crew of RV *Polarstern* (Knust, 2017) and participants of leg one to three for their help in the field. We would like to thank Joshua M. L. King for insightful discussions and comments.



## References

- Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S., Davis, A., Dean, J., Devin, M., Ghemawat, S., Goodfellow, I., Harp, A., Irving, G., Isard, M., Jia, Y., Jozefowicz, R., Kaiser, L., Kudlur, M., Levenberg, J., Mané, D., Monga, R., Moore, S., Murray, D., Olah, C., Schuster, M., Shlens, J., Steiner, B., Sutskever, I., Talwar, K., Tucker, P., Vanhoucke, V., Vasudevan, V., Viégas, F., Vinyals, O., Warden, P., Wattenberg, M., Wicke, M., Yu, Y., and Zheng, X.: TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems, <http://tensorflow.org/>, software available from tensorflow.org, 2015.
- 435 Bahdanau, D., Cho, K., and Bengio, Y.: Neural machine translation by jointly learning to align and translate, arXiv preprint arXiv:1409.0473, 2014.
- Bengio, Y., Delalleau, O., and Le Roux, N.: 11 label propagation and quadratic criterion, 2006.
- 440 Bishop, C. M.: Pattern recognition and machine learning, Springer, 2006.
- Breiman, L.: Random forests, *Machine learning*, 45, 5–32, 2001.
- Calonne, N., Richter, B., Löwe, H., Cetti, C., ter Schure, J., Van Herwijnen, A., Fierz, C., Jaggi, M., and Schneebeli, M.: The RHOSSA campaign: multi-resolution monitoring of the seasonal evolution of the structure and mechanical stability of an alpine snowpack, *The Cryosphere*, 14, 1829–1848, 2020.
- 445 Chen, C., Liaw, A., Breiman, L., et al.: Using random forest to learn imbalanced data, University of California, Berkeley, 110, 24, 2004.
- Chollet, F. et al.: Keras, <https://github.com/fchollet/keras>, 2015.
- Colbeck, S.: A review of the metamorphism and classification of seasonal snow cover crystals, IAHS Publication, 162, 3–24, 1987.
- Cortes, C. and Vapnik, V.: Support-vector networks, *Machine learning*, 20, 273–297, 1995.
- Cover, T. and Hart, P.: Nearest neighbor pattern classification, *IEEE Transactions on Information Theory*, 13, 21–27,
- 450 <https://doi.org/10.1109/TIT.1967.1053964>, 1967.
- CyberZHG: Keras Self-Attention, <https://github.com/CyberZHG/keras-self-attention>, 2020.
- Douville, H., Royer, J.-F., and Mahfouf, J.-F.: A new snow parameterization for the Meteo-France climate model, *Climate Dynamics*, 12, 21–35, 1995.
- Fierz, C., Armstrong, R. L., Durand, Y., Etchevers, P., Greene, E., McClung, D. M., Nishimura, K., Satyawali, P. K., and Sokratov, S. A.: The international classification for seasonal snow on the ground, 2009.
- 455 Fix, E. and Hodges Jr, J. L.: Discriminatory analysis-nonparametric discrimination: Small sample performance, Tech. rep., CALIFORNIA UNIV BERKELEY, 1952.
- Forgy, E. W.: Cluster analysis of multivariate data: efficiency versus interpretability of classifications, *biometrics*, 21, 768–769, 1965.
- Ghahramani, Z.: *Unsupervised Learning*, pp. 72–112, Springer Berlin Heidelberg, Berlin, Heidelberg, [https://doi.org/10.1007/978-3-540-28650-9\\_5](https://doi.org/10.1007/978-3-540-28650-9_5), 2004.
- 460 Han, J., Kamber, M., and Pei, J.: 9 - Classification: Advanced Methods, in: *Data Mining (Third Edition)*, edited by Han, J., Kamber, M., and Pei, J., The Morgan Kaufmann Series in Data Management Systems, pp. 393–442, Morgan Kaufmann, Boston, third edition edn., <https://doi.org/https://doi.org/10.1016/B978-0-12-381479-1.00009-5>, 2012.
- Havens, S., Marshall, H.-P., Steiner, N., and Tedesco, M.: Snow micro penetrometer and near infrared photography for grain type classification, in: *2010 International Snow Science Workshop*, pp. 465–469, 2010.
- 465 Havens, S., Marshall, H.-P., Pielmeier, C., and Elder, K.: Automatic grain type classification of snow micro penetrometer signals with random forests, *IEEE transactions on geoscience and remote sensing*, 51, 3328–3335, 2012.



- Herla, F., Horton, S., Mair, P., and Haegeli, P.: Snow profile alignment and similarity assessment for aggregating, clustering, and evaluating snowpack model output for avalanche forecasting, *Geoscientific Model Development*, 14, 239–258, [https://doi.org/10.5194/gmd-14-239-](https://doi.org/10.5194/gmd-14-239-2021)  
470 2021, 2021.
- Ho, T. K.: Random decision forests, in: *Proceedings of 3rd international conference on document analysis and recognition*, vol. 1, pp. 278–282, IEEE, 1995.
- Hochreiter, S. and Schmidhuber, J.: Long short-term memory, *Neural computation*, 9, 1735–1780, 1997.
- Ismail Fawaz, H., Forestier, G., Weber, J., Idoumghar, L., and Muller, P.-A.: Deep learning for time series classification: a review, *Data mining and knowledge discovery*, 33, 917–963, 2019.
- Johnson, J. B. and Schneebeli, M.: Snow strength penetrometer, uS Patent 5,831,161, 1998.
- Jurafsky, D. and Martin, J. H.: *Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition*, <https://web.stanford.edu/~jurafsky/slp3/>, in progress. 3rd ed. draft. Can be found at <https://web.stanford.edu/~jurafsky/slp3/>, 2021.
- 480 Kaltenborn, J., Macfarlane, A. R., Clay, V., and Schneebeli, M.: Pre-trained Models for SMP Classification and Segmentation, <https://doi.org/10.5281/zenodo.7063521>, 2022.
- King, J., Howell, S., Brady, M., Toose, P., Derksen, C., Haas, C., and Beckers, J.: Local-scale variability of snow density on Arctic sea ice, *The Cryosphere*, 14, 4323–4339, 2020.
- Knust, R.: Polar research and supply vessel POLARSTERN operated by the Alfred-Wegener-Institute, *Journal of large-scale research facilities JLSRF*, 3, A119–A119, 2017.
- 485 Lemaître, G., Nogueira, F., and Aridas, C. K.: Imbalanced-learn: A python toolbox to tackle the curse of imbalanced datasets in machine learning, *The Journal of Machine Learning Research*, 18, 559–563, 2017.
- Li, D., Hasanaj, E., and Li, S.: 3 – Baselines, <https://blog.ml.cmu.edu/2020/08/31/3-baselines/>, [Online; <https://blog.ml.cmu.edu/2020/08/31/3-baselines/>, accessed 04-March-2021], 2020.
- 490 Liu, X.-Y., Wu, J., and Zhou, Z.-H.: Exploratory undersampling for class-imbalance learning, *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)*, 39, 539–550, 2008.
- Lloyd, S.: Least squares quantization in PCM, *IEEE transactions on information theory*, 28, 129–137, 1982.
- Löwe, H. and Van Herwijnen, A.: A Poisson shot noise model for micro-penetration of snow, *Cold Regions Science and Technology*, 70, 62–70, 2012.
- 495 Macfarlane, A. R., Schneebeli, M., Dadic, R., Wagner, D. N., Arndt, S., Clemens-Sewall, D., Hämmerle, S., Hannula, H.-R., Jaggi, M., Kolabutin, N., Krampe, D., Lehning, M., Matero, I., Nicolaus, M., Oggier, M., Pirazzini, R., Polashenski, C., Raphael, I., Regnery, J., Shimanchuck, E., Smith, M. M., and Tavri, A.: Snowpit SnowMicroPen (SMP) force profiles collected during the MOSAiC expedition, PANGAEA, <https://doi.org/10.1594/PANGAEA.935554>, in: Macfarlane, AR et al. (2021): Snowpit raw data collected during the MOSAiC expedition. PANGAEA, <https://doi.org/10.1594/PANGAEA.935934>, 2021.
- 500 Ménard, C. B., Essery, R., Barr, A., Bartlett, P., Derry, J., Dumont, M., Fierz, C., Kim, H., Kontu, A., Lejeune, Y., et al.: Meteorological and evaluation datasets for snow modelling at 10 reference sites: description of in situ and bias-corrected reanalysis data, *Earth System Science Data*, 11, 865–880, 2019.
- Nguyen, N. and Guo, Y.: Comparisons of sequence labeling algorithms and extensions, in: *Proceedings of the 24th international conference on Machine learning*, pp. 681–688, 2007.



- 505 Nicolaus, M., Perovich, D. K., Spreen, G., Granskog, M. A., von Albedyll, L., Angelopoulos, M., Anhaus, P., Arndt, S., Belter, H. J., Bessonov, V., Birnbaum, G., Brauchle, J., Calmer, R., Cardellach, E., Cheng, B., Clemens-Sewall, D., Dadic, R., Damm, E., de Boer, G., Demir, O., Dethloff, K., Divine, D. V., Fong, A. A., Fons, S., Frey, M. M., Fuchs, N., Gabarró, C., Gerland, S., Goessling, H. F., Gradinger, R., Haapala, J., Haas, C., Hamilton, J., Hannula, H.-R., Hendricks, S., Herber, A., Heuzé, C., Hoppmann, M., Høyland, K. V., Huntemann, M., Hutchings, J. K., Hwang, B., Itkin, P., Jacobi, H.-W., Jaggi, M., Jutila, A., Kaleschke, L., Katlein, C., Kolabutin, N., Krampe, D.,
- 510 Kristensen, S. S., Krumpfen, T., Kurtz, N., Lampert, A., Lange, B. A., Lei, R., Light, B., Linhardt, F., Liston, G. E., Loose, B., Macfarlane, A. R., Mahmud, M., Matero, I. O., Maus, S., Morgenstern, A., Naderpour, R., Nandan, V., Niubom, A., Oggier, M., Oppelt, N., Pätzold, F., Perron, C., Petrovsky, T., Pirazzini, R., Polashenski, C., Rabe, B., Raphael, I. A., Regnery, J., Rex, M., Ricker, R., Riemann-Campe, K., Rinke, A., Rohde, J., Salganik, E., Scharien, R. K., Schiller, M., Schneebeil, M., Semmling, M., Shimanchuk, E., Shupe, M. D., Smith, M. M., Smolyanitsky, V., Sokolov, V., Stanton, T., Stroeve, J., Thielke, L., Timofeeva, A., Tonboe, R. T., Tavri, A., Tsamados, M., Wagner, D. N., Watkins, D., Webster, M., and Wendisch, M.: Overview of the MOSAiC expedition: Snow and sea ice, *Elementa: Science of the Anthropocene*, 10, <https://doi.org/10.1525/elementa.2021.000046>, 000046, 2022.
- Nixdorf, U., Dethloff, K., Rex, M., Shupe, M., Sommerfeld, A., Perovich, D. K., Nicolaus, M., Heuzé, C., Rabe, B., Loose, B., Damm, E., Gradinger, R., Fong, A., Maslowski, W., Rinke, A., Kwok, R., Spreen, G., Wendisch, M., Herber, A., Hirsekorn, M., Mohaupt, V., Frickenhaus, S., Immerz, A., Weiss-Tuider, K., König, B., Mengedoht, D., Regnery, J., Gerchow, P., Ransby, D., Krumpfen, T., Morgenstern, A., Haas, C., Kanzow, T., Rack, F. R., Saitzev, V., Sokolov, V., Makarov, A., Schwarze, S., Wunderlich, T., Wurr, K., and Boetius, A.: MOSAiC Extended Acknowledgement, <https://doi.org/10.5281/zenodo.5541624>, 2021.
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., and Duchesnay, E.: Scikit-learn: Machine Learning in Python, *Journal of Machine Learning Research*, 12, 2825–2830, 2011.
- 525 Pörtner, H.-O., Roberts, D. C., Masson-Delmotte, V., Zhai, P., Tignor, M., Poloczanska, E., and Weyer, N.: The ocean and cryosphere in a changing climate, *IPCC Special Report on the Ocean and Cryosphere in a Changing Climate*, 2019.
- Russell, S. and Norvig, P.: *Artificial intelligence: a modern approach*, 2002.
- Satyawali, P., Schneebeil, M., Pielmeier, C., Stucki, T., and Singh, A.: Preliminary characterization of Alpine snow using SnowMicroPen, *Cold Regions Science and Technology*, 55, 311–320, 2009.
- 530 Schneebeil, M. and Johnson, J. B.: A constant-speed penetrometer for high-resolution snow stratigraphy, *Annals of Glaciology*, 26, 107–111, 1998.
- Schneebeil, M., Pielmeier, C., and Johnson, J. B.: Measuring snow microstructure and hardness using a high resolution penetrometer, *Cold Regions Science and Technology*, 30, 101–114, 1999.
- Schölkopf, B., Smola, A. J., Bach, F., et al.: *Learning with kernels: support vector machines, regularization, optimization, and beyond*, MIT press, 2002.
- 535 Schuster, M. and Paliwal, K. K.: Bidirectional recurrent neural networks, *IEEE transactions on Signal Processing*, 45, 2673–2681, 1997.
- Soni, R. and Mathai, K. J.: Improved Twitter sentiment prediction through cluster-then-predict model, *arXiv preprint arXiv:1509.02437*, 2015.
- Steger, C., Kotlarski, S., Jonas, T., and Schär, C.: Alpine snow cover in a changing climate: a regional climate model perspective, *Climate dynamics*, 41, 735–754, 2013.
- 540 Stone, M.: Cross-validated choice and assessment of statistical predictions, *Journal of the Royal Statistical Society: Series B (Methodological)*, 36, 111–133, 1974.





- Sturm, M. and Massom, R. A.: Snow in the sea ice system: Friend or foe, *Sea ice*, pp. 65–109, 2017.
- Theodorou, T., Mporas, I., and Fakotakis, N.: An overview of automatic audio segmentation, *International Journal of Information Technology and Computer Science (IJITCS)*, 6, 1, 2014.
- 545
- Trivedi, S., Pardos, Z. A., and Heffernan, N. T.: The utility of clustering in prediction tasks, arXiv preprint arXiv:1509.06163, 2015.
- Wu, X., Kumar, V., Quinlan, J. R., Ghosh, J., Yang, Q., Motoda, H., McLachlan, G. J., Ng, A., Liu, B., Philip, S. Y., et al.: Top 10 algorithms in data mining, *Knowledge and information systems*, 14, 1–37, 2008.
- Yarowsky, D.: Unsupervised word sense disambiguation rivaling supervised methods, in: 33rd annual meeting of the association for computational linguistics, pp. 189–196, 1995.
- 550
- Zhou, D., Bousquet, O., Lal, T. N., Weston, J., and Schölkopf, B.: Learning with local and global consistency, in: *Advances in Neural Information Processing Systems 16*, pp. 321–328, MIT Press, 2004.
- Zhu, X. J. and Ghahramani, Z.: Learning from labeled and unlabeled data with label propagation, 2002.