# Automated Prediction of Gamburtsev Subglacial Lakes in East Antarctica With Optimized Stacking Ensemble Learning

Qian Ma, Tong Hao, *Member, IEEE*, Tiantian Feng, Gang Qiao, *Member, IEEE*, Asoke K. Nandi, *Life Fellow, IEEE*, and Chen Lv

*Abstract*— The development in machine learning (ML) technology has brought new horizons for the prediction of subglacial lakes (SLs) using radio-echo sounding (RES) data, offering fresh perspectives toward the automated identification of SLs. Nonetheless, the inherent data imbalance across various classes within the dataset presents significant analytical challenges. To address this limitation, the artificial bee colony (ABC) optimization algorithm is introduced to automatically predict SLs in Gamburtsev Province in East Antarctica, using an optimized stacking ensemble learning approach. The proposed method predicts SLs by using five representative features selected through importance and correlation analyses of eight features derived from RES data. The experimental outcomes demonstrate the superiority of this method in overcoming the significant imbalance of RES data, successfully identifying known lakes in the validation dataset. Furthermore, this study summarizes an inventory of SLs across the Gamburtsev subglacial mountains in East Antarctica, and a total of 55 new candidate SLs with lengths ranging from 108 to 38 130 m have been predicted using our novel method. The source code is publicly available at https://github.com/vivian-ma97/ABC-Stacking-for-Subglacial-Lakes

*Index Terms*— Antarctica, ensemble learning, feature extraction, radio echo sounding, subglacial lakes.

## I. INTRODUCTION

SINCE the discovery of the first subglacial lake (SL) in Antarctica, the prediction and identification of subglacial water systems have attracted worldwide attention from geophysicists, glaciologists, microbiologists, and geologists [1],

Qian Ma and Tong Hao are with Shanghai Research Institute for Intelligent Autonomous Systems, Center for Spatial Information Science and Sustainable Development Applications, College of Surveying and Geo-Informatics, Tongji University, Shanghai 200092, China (e-mail: qianma@tongji.edu.cn; tonghao@tongji.edu.cn).

Tiantian Feng, Gang Qiao, and Chen Lv are with the Center for Spatial Information Science and Sustainable Development Applications, College of Surveying and Geo-Informatics, Tongji University, Shanghai 200092, China (e-mail: fengtiantian@tongji.edu.cn; qiaogang@tongji.edu.cn; 2231970@tongji.edu.cn).

Asoke K. Nandi is with the College of Engineering, Design and Physical Sciences, Brunel University of London, UB8 3PH Uxbridge, U.K. (e-mail: Asoke.Nandi@brunel.ac.uk).

Digital Object Identifier 10.1109/TGRS.2025.3587133

[2], [3], [4], with the motivation stemming from their importance in glaciology, biology, climatology, and geology. For instance, SLs are one of the key integral components of the subglacial hydrological system, exerting significant influence on the dynamics, evolution, flow velocity, and mass balance of ice sheets [5].

The formation of SLs is attributed to a confluence of multifaceted factors, such as the subglacial topography, geothermal heat flux, insulation, and pressure of the overlying ice sheet [6], [7], [8]. For instance, subglacial topographic basins offer natural gathering corners of subglacial waters that may be melted from ice by geothermal heat transferred from the Earth's interior, and it is such heat that controls the ice rheology and hence decouples the ice sheet from the bedrock in the context of subglacial waters [9].

Remote sensing and geophysical technologies have revolutionized the observation of the cryosphere, with advanced equipments being developed and employed on a wide range of ground-based, airborne, and orbital platforms [10], [11], [12]. Among these, radio-echo sounding (RES) is widely recognized as an effective tool to shed light on the hard-to-observe subglacial water-saturated environment below kilometers of ice sheet [13]. Since the first inventory of 17 SLs in East Antarctica discovered by RES surveys [14], there have been five comprehensive SLs inventories reported to date [13], [14], [15], [16], [17]. In the most recent global inventory, there are 675 possible SLs located in Antarctica. More than 80% of the stable lakes, implying either a closed hydraulic system or approximately balanced inflows and outflows, have been predominantly detected by RES [13].

Although the widely adopted visual inspection and interpretation method offers an intuitive means to identify SLs by establishing some predefined picking criteria, such as relative brightness, topographic/hydraulic flatness, and echo abruptness [18], this approach is subjective and time-consuming, making itself impractical for modern large-scale RES data collection [19], [20]. Consequently, more sophisticated semi-automated or fully automated methodologies have been proposed, implemented, and optimized to enhance the efficiency and accuracy of SL prediction from large-volume RES B-scope data. For instance, Carter et al. [21] developed an automated technique leveraging the distinctive radar reflection characteristics of different basal materials for a systematical categorization of SLs. Among various parameters, the basal

radar reflectivity is the most important quantitative metric for the prediction of subglacial water [12]. However, the dielectric loss in the overlying ice layers complicates its accurate calculation. In 2013, Wolovick et al. [18] further refined the prediction process by combining reflectivity analysis with manual digitization. In their reflectivity analysis, highreflectivity anomalies were used to identify potential SLs in RES B-scope data. Simultaneously, their manual digitization approach involved an operator who not only examined the radargrams but also selected basal reflectors as potential water bodies based on their morphological alignment with predefined criteria. This work provided a comprehensive inventory of SLs detected in Antarctica's Gamburtsev Province (AGAP) region. Building upon this foundational work, Livingstone et al. [13] expanded to a worldwide SL inventory with a systematic three-stage identification method. In the first stage, they conducted a visual inspection on RES data, manually looking for continuous, bright, and flat reflectors—features commonly associated with SLs. In the second stage, they evaluated potential lake candidates by analyzing bed return power (BRP) in relation to surrounding topographic variations. Finally, in the third stage, they examined the BRP variability to assess the SL candidates. On the other hand, for distinct characteristics of the ice–water interface compared with the ice–bedrock interface, numerous studies have endeavored to identify SLs by synergistically combining basal topographic features [22], [23], RES waveform shapes [24], [25], reflection characteristics [26], [27], and roughness [28]. Despite the progress made with manual and semi-automated methods for SLs prediction, the limitations of these approaches become evident when faced with vast datasets typical of modern RES data collection. This has paved the way for the development of more advanced technologies that leverage both spatial and temporal data characteristics.

In contrast to traditional SL prediction methods, which solely rely on statistical information from the 2-D RES B-scope data, such as BRP, reflectivity, roughness, and hydraulic head, the inherent time–frequency characteristics of 1-D RES signals (A-scope) are often overlooked. Starting from this, we recently proposed an automated method for predicting SLs utilizing joint time–frequency analysis (JTFA) methods, in particular, the short-time Fourier transform (STFT) technique [29]. The STFT is a powerful mathematical tool that allows for the analysis of time-varying signals in the frequency domain [30]. By applying the STFT to the RES data, we can extract the frequency characteristics and temporal variations associated with the SLs, which is effectively a simpler quantification of the A-scope shape feature utilized in [31]. In our previous work [29], we defined STFT feature with an empirically determined variable which also integrated the time–frequency information with terrain slope. Although such an empirical formula of the STFT feature successfully demonstrated good accuracy in subglacial water bodies, the method itself lacks intuitive physical meaning and relies on selected empirical thresholds for decision-making. Therefore, it potentially can lead to the overestimation or underestimation of subglacial water bodies, thereby affecting the prediction accuracy. These limitations spark the motivation to explore new methods to automatically optimize the prediction process with reduced reliance on empirical parameters and enhanced objectivity and adaptability of the model.

On the other hand, machine learning (ML) techniques are becoming another popular stream of investigation on SLs [22], [31], [32]. For instance, Ilisei et al. [31] developed an automated SL prediction technique employing ML techniques, in particular, the support vector machine (SVM) algorithm. Their methods aim to enhance the prediction of SLs by automatically extracting a defined set of discriminant features from the ice–basal interface, including terrain, shape, and statistical characteristics. These features are used to capture the unique signatures of SLs in RES data, which the SVM algorithm can classify. Their automatic approach ensures consistent results for RES datasets across two Antarctic regions. The ratio of positive and negative samples in their training sample is about 1:2. Lang et al. [22] proposed a semiautomatic approach based on local and regional characteristics to identify the subglacial dry–wet transition zones. By combining an SVM classifier with the synthetic minority over-sampling technique, they successfully demonstrated the potential of ML models in dealing with the issue of data imbalance. In their case, the classification effect is verified when the ratio of positive and negative samples is approximately 1:3. More recently, Dong et al. [32] identified the distribution of different categories of the basal reflector features using an unsupervised encode-cluster algorithm, effectively mitigating the imbalance between SLs and non-SLs. Unsupervised learning methods are generally considered to have strong generalization capabilities; however, some researchers claim that unsupervised learning methods might fall short of supervised learning methods in terms of precise prediction [33]. Therefore, how to achieve reliable classification results for SLs using such an extremely imbalanced number of positive and negative samples remains an area for further study.

The stacking method, one of the representative methods of ensemble learning, can construct and integrate multiple learners through a predefined combination strategy [34]. If these base learners are well-differentiated, the overall model will be more robust to data noise overfitted by any single learner [35]. This efficacy stems from their capacity to amalgamate diverse and heterogeneous algorithms, creating a synergistic effect that enhances predictive accuracy. In addition, the stacking model combines the predictions of multiple different base learners to form the final output, which can alleviate the problem of imbalanced classification of samples. This is because different base learners may have different sensitivities to different majority or minority parts of the dataset. Especially when dealing with imbalanced data, by combining the predictions of these base learners, the stacking model is able to capture the characteristics of the dataset fully [36].

In this study, the artificial bee colony (ABC) optimization algorithm is adopted to automatically predict SLs using an optimized stacking ensemble learning approach, which effectively handles extremely imbalanced positive-to-negative sample ratios. First, five representative features for SL prediction are identified through importance and correlation analyses of eight features derived from RES data.

Subsequently, the ABC optimization algorithm is employed to select the optimum ensemble learning model, as well as the most suitable base learner and metalearner for the following stacking fusion, which is then utilized for SL prediction. To validate the effectiveness of our proposed method, we conducted a set of experiments and evaluated the results from both quantitative and qualitative perspectives. The experimental results demonstrate that the optimization algorithm can effectively identify SLs, demonstrating strong stability even in cases of highly imbalanced data distribution.

The main contributions of this study are as follows.

1) The proposed approach signifies a strategic shift in addressing the issue of data imbalance, focusing on model architecture adjustments to accommodate and mitigate the challenges posed by uneven data distributions.

2) Compared with existing ML approaches, we propose a stacking-based ensemble learning approach that aids in selecting the optimum classifier combination, contributing to an enlarged overall prediction accuracy even under conditions of extreme data imbalance.

3) An updated inventory of SLs in AGAP region, with an additional 55 new candidate SLs, is provided in this article.

The remainder of this article is organized as follows. Sections II and III present a comprehensive review of the characteristics of the basal interface, as well as the feature extraction and feature selection operations. Section IV provides a detailed description of the proposed stacking-based ensemble learning method for SL prediction. The experimental results on real RES data and subsequent discussions are presented in Section V. Finally, we summarize the key findings of this study in Section VI.

## II. CHARACTERISTICS OF THE BASAL INTERFACE

### A. Basal Topography and Hydraulic Characteristics

As illustrated in Fig. 1(a), SLs are primarily located in topographic depressions. In addition, the gradient of hydraulic heads in SLs is typically low, allowing water to accumulate [37]. This pattern of water accumulation, governed by the contours of the subglacial terrain and hydraulic head, promotes the formation of SLs in regions characterized by hydraulic flatness, which is a fundamental hydrological attribute of SLs. Roughness is a critical statistical parameter for describing the irregularity and complexity of subglacial topography [38]. Areas with lower roughness typically indicate smoother and more continuous terrain, which can facilitate the convergence of subglacial water flows and the formation of SLs. Conversely, areas with higher roughness are characterized by greater terrain undulations, which may impede the continuous flow of water, thus hindering lake formation.

### B. Basal Reflectivity Contrast

When a radar wave encounters the interface between ice, water, and bedrock, the permittivity contrast of these different materials leads to changes in the returned power of the echo signal. Assuming the vertical incidence of radar waves, the reflectivity $\rho$ at the interface between ice and basal material is expressed as

$$\rho = \left( \frac{\sqrt{\epsilon_{r_2}} - \sqrt{\epsilon_{r_1}}}{\sqrt{\epsilon_{r_2}} + \sqrt{\epsilon_{r_1}}} \right)^2 \tag{1}$$

where $\epsilon_{r_1}$ is the relative permittivity of ice and $\epsilon_{r_2}$ is the relative permittivity of the basal material, e.g., water or bedrock. Water provides a much sharper electromagnetic boundary due to its higher relative permittivity ($\sim$81), compared with bedrock (3$\sim$30) and ice ($\sim$3.17) [39]. It is worth noting that the dielectric permittivity of these subglacial materials may slightly differ from the values presented here, as they can be influenced by factors such as temperature fluctuations and material composition. For instance, the dielectric permittivity of water may increase to around 88 near the freezing point [40].

The difference in dielectric properties between ice and water/bedrock results in a change in the reflectivity and penetration of radar waves across the interface [27], as illustrated in (1). The high permittivity contrast can result in stronger reflected signals with more concentrated frequency components. In addition, specular reflections occur more often at the ice–water interface, while deeper penetration and multiple reflections are common for bedrocks [41]. When incident waves encounter discontinuities at rugged basal boundaries, coherent cancellation may occur, leading to a gradual weakening of the reflected signal strength, leading to a wider A-scope waveform, as illustrated in Fig. 1(b) [31].

## III. FEATURE EXTRACTION AND FEATURE SELECTION

In this section, we utilize a general feature selection procedure on the selected (sub)glacial features, with which further features for other applications can also be included, examined, and selected.

### A. Feature Extraction

Eight features encapsulating the most salient characteristics of SLs are included in this study: hydraulic head, hydraulic head gradient, topographic roughness, BRP, corrected BRP (CBRP), time–frequency characteristics of A-scope, ice thickness, and basal elevation. Here, we define the basal elevation as the elevation at the bottom of the overlying ice sheet, so it could be the elevation at the ice–water interface or the ice–bedrock interface.

*1) Hydraulic Head and Hydraulic Head Gradient:* Hydraulic head serves as the fundamental driving force for water movement and profoundly influences its direction [42], [43]. Subglacial water naturally flows from regions with higher hydraulic heads to those with lower ones. The areas characterized by a high hydraulic head are less likely to harbor water, whereas areas with a low subglacial hydraulic head are more prone to contain water. Consequently, the computation of subglacial hydraulic head can serve as an indicator of the presence possibility of subglacial water. Following the traditional method in [21] and [37], we use the hydraulic head and its gradient as features of SLs. The hydraulic head $\Phi$ at
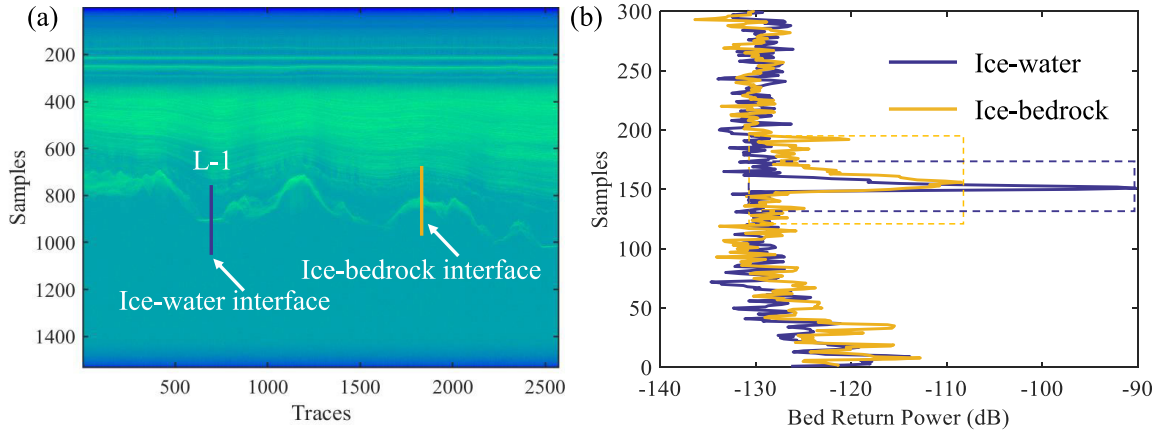
Fig. 1. (a) Typical B-scope radargram with one clear SL (L-1) in the Gamburtsev subglacial mountains in East Antarctica. (b) Comparison of two A-scope waveforms crossing the ice–water and ice–bedrock interfaces.

the bottom of an ice sheet is defined by its surface elevation ($S$) and the basal elevation beneath the ice ($B$), as

$$\phi = \frac{\rho_i}{\rho_w} S + \left(1 - \frac{\rho_i}{\rho_w}\right) B \qquad (2)$$

where the density of ice $\rho_i$ is 917 kg/m$^3$ and the density of water $\rho_w$ is 1000 kg/m$^3$.

The gradient of the hydraulic head reflects the direction of the water flow, and a smaller gradient means a smaller direction change in the hydraulic head. In this study, the gradient of the hydraulic head $\omega$ is defined as

$$\omega = \frac{\Delta z}{\Delta x} \qquad (3)$$

where $\Delta z$ is the difference in hydraulic head between two adjacent A-scope peaks, $\Delta x$ is the horizontal distance between two adjacent A-scopes, and the unit for both variables is in meters.

2) *Topographic Roughness:* The ice–water interface is normally relatively flat with a low topographic roughness in the SL area, while the ice–bedrock interface commonly has a higher topographic roughness. Here, we define the local elevation changes of the basal terrain as the roughness feature [37]. The local elevation changes $\zeta$ is calculated by

$$\zeta = \frac{1}{n-1} \sum_{i=1}^{n} [z(x_i) - \overline{z}] \qquad (4)$$

where $n$ is the number of A-scopes corresponding to the minimum predictable lake length. A survey of the existing lakes in the AGAP area inventory [13], [18] revealed that these lakes are all over 100 m in length. Therefore, if the horizontal resolution of the profile is 18 m, $n = 5$. If the horizontal resolution is 30 m, $n = 3$. $\overline{z}$ is the mean value of $z(x_i)$, which can be defined by

$$\overline{z}(x_i) = E(z(x_i)) \qquad (5)$$

where $E$ is the expectation operation.

3) *BRP and CBRP:* The BRP from RES surveys provides insights into the properties of the ice–basal interface, encompassing the basal material and the roughness of the interface [44]. Particularly strong changes in relative echo strength are commonly associated with changes in bed wetness, enabling radar techniques to be used to map SLs in Antarctica [45]. Considering the dissipation of radar waves within bedrock and the less specular reflections at the ice–bedrock interface, such complex interactions altogether lead to progressive attenuation of radar echo intensity, as can be discerned via the RES A-scope data [31]. In addition, adhering to the principles outlined in [41], CBRP is commonly implemented to compensate for the depth-related attenuation difference, which is formulated as

$$[P_a]_{\mathrm{dB}} = [P]_{\mathrm{dB}} + 2\left[\left(h + \frac{d}{\sqrt{\epsilon_{r_1}}}\right)\right]_{\mathrm{dB}} + 2d\langle N \rangle \qquad (6)$$

where $[P]_{\mathrm{dB}} = 10\log_{10}(P)$ denotes the recorded BRP in dB, $h$ represents the aircraft's height above the ice surface, $d$ is the thickness of the ice, $\epsilon_{r_1}$ is the relative permittivity of ice, and $N$ is 11.7 dB/km as calculated by Wolovick et al. [18] for the AGAP regional attenuation rate.

4) *Time–Frequency Characteristics:* To analyze the inherent frequency characteristics of each A-scope, we employed the time–frequency feature (TFF) derived from the STFT to quantify the shape of the A-scope waveform. Similar to Hao et al. [29], we utilized the time–frequency characteristics at the ice–bed interface. However, unlike the empirical equation [29, eq. (9)] in [29], we defined the TFF in (7) with two key parameters, namely, $F$, the normalized maximum frequency of the A-scope, and $A$, the magnitude of that maximum frequency. By doing so, we circumvented the use of empirical parameters established by Hao et al. [29], offering a more direct and objective approach for analyzing the time–frequency characteristics

$$\mathrm{TFF} = F \times A. \qquad (7)$$

Fig. 2 elucidates the procedure of calculating the STFT response. The waveform at the ice–water and ice–bedrock interfaces is presented in Fig. 2(a) and (d). To ensure that
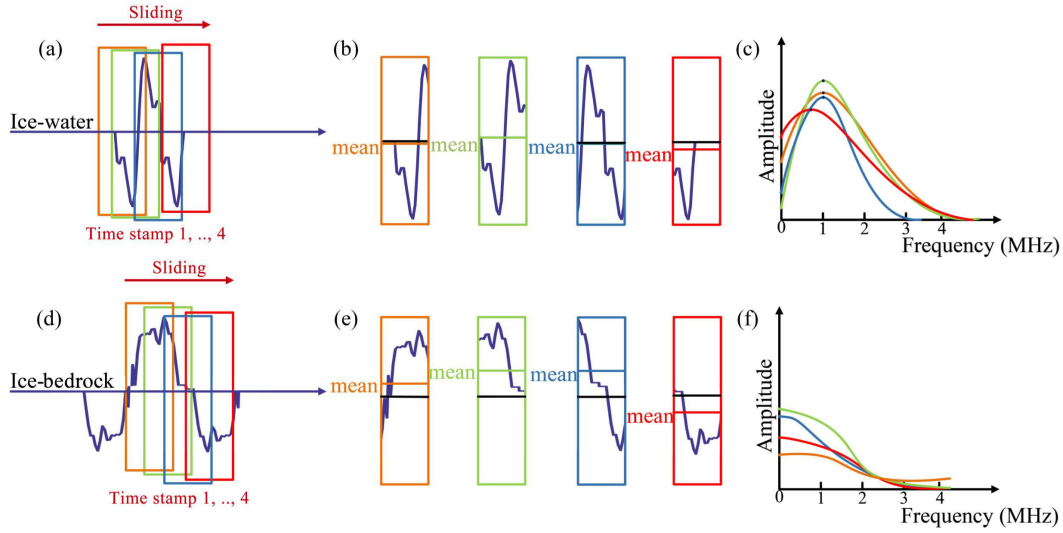
Fig. 2. Underlying mechanism of time–frequency characteristics of ice–water and ice–bedrock interfaces. Four differently positioned FFT windows to section the A-scope waveform for (a) ice–water and (d) ice–bedrock interfaces. The corresponding sectioned waveforms and mean values for (b) ice–water and (e) ice–bedrock interfaces. The corresponding frequency responses for (c) ice–water and (f) ice–bedrock interfaces.
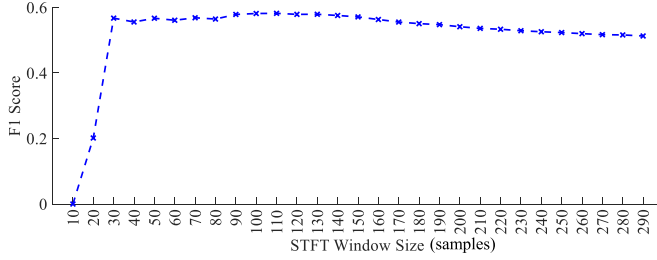


Fig. 3. $F1$ scores calculated using different STFT window sizes.

the mean value within each time window is close to 0, the sliding window, represented by different colors, is placed around the main peak of the A-scope waveform at different time steps along the range direction. In this study, we chose to utilize a Hanning window due to its ability to effectively reduce high-frequency interference and spectral leakage. Fig. 2(b) and (e) displays the mean values of the signal at these time points, where the decrease in mean values indicates significant changes in signal strength or characteristics over time. Fig. 2(c) illustrates the frequency response at the ice–water interface, where the highest response no longer appears at zero frequency. This is attributed to the presence of strong high-frequency components in the ice–water interface signal, resulting in an enhanced response at nonzero frequencies. In contrast, as shown in Fig. 2(d), the signal behavior at the ice–bedrock interface is different. Here, the window either fails to fully cover the entire main peak as shown in Fig. 2(e), or the mean value within the window is quite high, resulting in a very strong zero-frequency component as depicted in Fig. 2(f). Therefore, the STFT response at the zero-frequency component can be used to distinguish between the ice–water interface and the ice–bedrock interface.

In addition, we investigate the effect of window size on capturing the signal characteristics of the ice–water and ice–bedrock interfaces. By comparing the known SL inventory

with the identified results using the TFF, we evaluate the correspondence between the TFF results and the labels for different window sizes. As shown in Fig. 3, the $F1$ score stabilizes and remains nearly unchanged when the window size is set to 30 samples or larger values, indicating that this window size of 30 samples effectively captures the signal characteristics. Further increasing the window size has minimal impact on model performance.

*5) Basal Elevation and Ice Thickness:* The basal elevation and ice thickness are features integral to the overall analysis because they may considerably influence the attenuation of the RES signals, which in turn has a significant impact on the prediction of subglacial waters. The ice thickness is determined using two two-way travel times of the radar signal: one for the radar wave traveling from the device to the ice surface and return and the other for the wave traveling from the device through the ice to the bottom and returning to the device.

By calculating the difference between these two timestamps, the total travel time of the radar waves through the entire ice layer is determined. This duration, when multiplied by the velocity of radar waves within the ice, allows for the estimation of the actual distance between the ice surface and the ice bottom, corresponding to the ice thickness. Basal elevation is then calculated by subtracting this derived ice thickness from the surface elevation.

### B. Feature Selection

In order to decrease the model's complexity and avoid the multicollinearity problem that redundant features might cause, it is necessary to first assess the importance and correlation of these features to identify the suitable input features for our model. As shown in Fig. 4, we initially used the Lasso regression for feature importance analysis. After determining the importance scores, we conducted a Pearson correlation analysis on eight features. As presented in Fig. 5, over 80% of
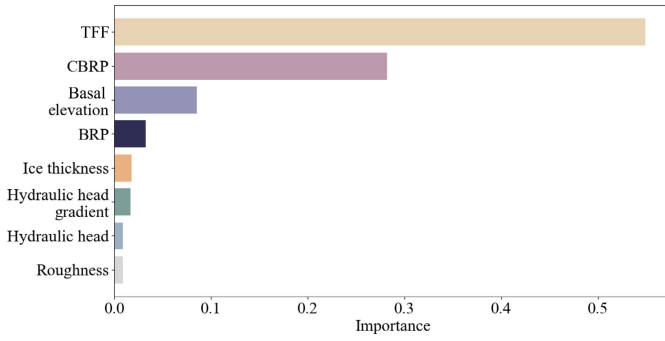
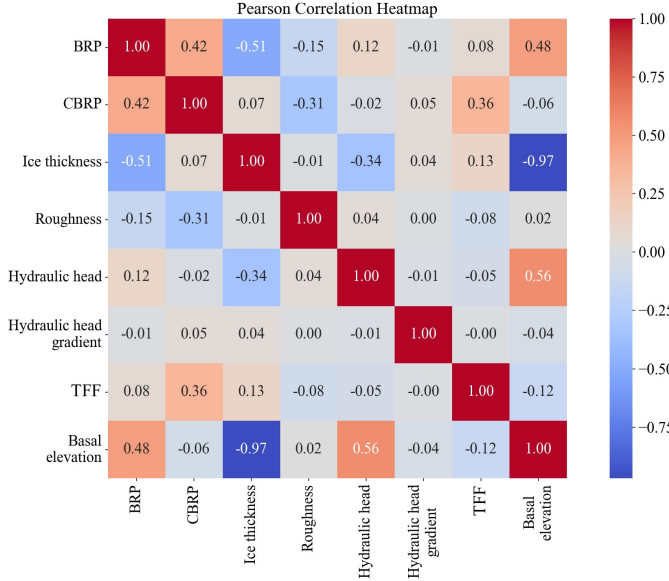Fig. 4.   Result of feature importance analysis by Lasso.



Fig. 5.   Result of Pearson correlation analysis.

the correlation values are below 0.4, with only five correlation coefficients being relatively high, specifically between BRP and CBRP, BRP and ice thickness, BRP and basal elevation, ice thickness and basal elevation, and hydraulic head and basal elevation. Combing with the feature importance presented in Fig. 4, we removed BPR, ice thickness, and hydraulic head. Finally, five features—TFF, CBRP, roughness, basal elevation, and hydraulic head gradient—are retained.

## IV. ABC-ASSISTED STACKING ENSEMBLE LEARNING METHOD

The stacking ensemble learning technique can enhance overall prediction performance by combining the outputs of multiple base models with a metalearner. This method is particularly advantageous for handling imbalanced data tasks, as it leverages the strengths of various models to improve the recognition of minority class samples more effectively [34], [46]. However, achieving effective collaboration among diverse base models and selecting the most complementary model combinations are key challenges in constructing a stacking model. To navigate the complexities in stacking and optimize its configuration, researchers have increasingly turned to metaheuristic optimization algorithms [47]. These algorithms

excel in exploring the vast combinatorial space of possible model configurations, thereby identifying the most potent combinations.

In line with this advanced approach, we employed the ABC algorithm [48], a metaphorical, honey bee-inspired metaheuristic optimization framework, to ascertain the optimal assemblage of base learners and suitable metalearners. By leveraging the ABC algorithm, we aim to harness its robust optimization capabilities to fine-tune our stacking ensemble, ensuring that we utilize the best possible combination of models and algorithms for our study's objectives. The optimization process for the ABC-assisted stacking ensemble learning method is shown in Fig. 6.

First, eight classic tree-based models—adaptive boosting (AdaBoost), categorical boosting (CatBoost), light gradient boosting machine (LightGBM), extreme gradient boosting (XGBoost), random forest (RF), extra trees (ETs), gradient boosting decision tree (GBDT), and decision tree (DT)— are selected as the initial configuration for the optimization process. The rationale behind this choice is that these models are well-known for their data augmentation capabilities, which is particularly advantageous in addressing issues of data imbalance [49], [50], [51].

Then, as shown in Fig. 6, we randomly select a configuration from the eight models and calculate its $F1$ score to assess its performance in the model selection part. Next, we optimize the model configurations using two different search strategies: local search and global search. During the local search phase, we first calculate the fitness value $F(x_i)$ of each classifier using (8), where the fitness value is calculated to represent the classifier's performance based on the $F1$ score. Specifically, the $F1$ score is used as an objective measure to assess the balance between precision and recall for each classifier, ensuring that their performance is evaluated. Then, the probability $P_i$ of each classifier being selected is calculated using (9). Those classifiers with higher $F1$ scores have a greater chance of being selected; thereby, these high-quality classifiers are prioritized for further local search

$$F(x_i) = \begin{cases} 1/(1 + f(x_i)) & \text{if } f(x_i) \geq 0 \\ 1 + |f(x_i)| & \text{if } f(x_i) < 0 \end{cases} \qquad (8)$$

where $F(x_i)$ represents the fitness value which can be calculated by $f(x_i)$, and $f(x_i)$ is the $F1$ score

$$P_i = \frac{F(x_i)}{\sum_{i=1}^{n} F(x_i)} \qquad (9)$$

where $n$ is the number of candidate classifiers currently in the optimization process.

At the same time, the global search seeks new possibilities by randomly selecting new model configurations. After each round of search, we calculate the $F1$ score of the new configuration and compare it with the original configuration, retaining the one that performs better. The search process is optimized by iteratively repeating the above steps. After each iteration, the algorithm evaluates the $F1$ scores of all model configurations and retains the configuration with the best performance. Furthermore, the optimal combination is obtained through iterative refinement, with ET, CatBoost, and
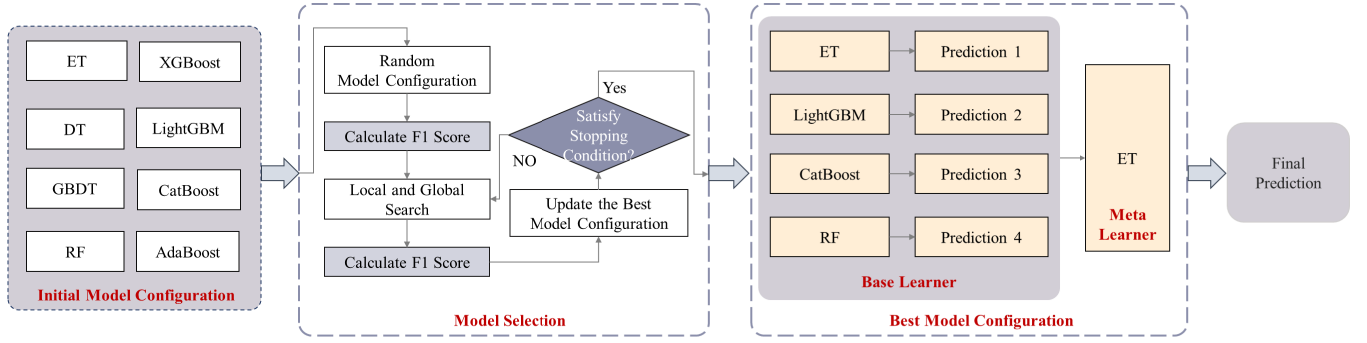
Fig. 6. Flowchart of the ABC-assisted stacking ensemble learning method.

LightGBM selected as the base learners, and ET serving as the meta learner. Finally, this set of models is applied to predict SLs in other regions of the AGAP.

## V. EXPERIMENTS AND RESULTS

To validate the performance of our proposed method, we conducted tests on RES datasets from the AGAP region in East Antarctica. This region is renowned for its complex subglacial topography, which includes numerous subglacial mountains and an extensive network of subglacial water systems [2]. The area has been extensively surveyed by airborne RES as part of the Gamburtsev Airborne Geophysical Mapping of Bedrock and Ice Targets (GAMBIT) Project, with longitudinal spacing of 33 km and lateral spacing of 5 km [18].

The primary data utilized in this study were sourced from the Multi-Channel Radar Depth Sounder (MCRDS), developed by the Center for Remote Sensing of Ice Sheets (CReSIS) at the University of Kansas [52], [53], [54], [55]. To provide a comprehensive summary of the SLs in the AGAP region, we conducted an in-depth analysis of the existing datasets. Although the data from Livingstone et al. [13] and Wolovick et al. [18] were collected during the same flight missions, they were processed using different focusing algorithms. Lamont-Doherty Earth Observatory (LDEO) employed the 1D-SAR algorithm, while CReSIS used the FK migration algorithm. To ensure consistency and comparability, we used the survey lines provided by CReSIS as a baseline and supplemented them with lines from the LDEO dataset that contains SLs not found in the CReSIS inventory. This effort specifically targeted those lines containing SLs identified by Wolovick et al. [18]. Through this process, we added a total of 18 B-scopes files and 25 517 A-scopes data entries. Overall, the dataset includes 923 RES B-scopes and contains a total of 1 573 006 A-scopes. A full description of the radar system's specifications and the RES dataset can be found in [54], which are also summarized in Table I.

In our model, we adopted all the SL locations provided by Livingstone et al. [13] and Wolovick et al. [18] as labels in our training and validation process. Specifically, only the survey lines containing inventoried SLs were selected, and each survey line was treated as a distinct dataset in a.mat file provided by CReSIS and LDEO, thus ensuring spatial independence. A 1:2 ratio was applied to randomly select a subset of these survey lines for the training and validation datasets,

TABLE I
PARAMETERS OF THE RES SYSTEM FOR THE AGAP REGION

| Parameter | Description |
|---|---|
| Campaign | AGAP |
| Radar system | MCRDS |
| Number of B-scope radargrams | 923 |
| Number of valid A-scopes | 1,573,006 |
| Platform type | Twin Otter aircraft |
| Platform height above ice sheet surface | Varied, hundreds of meters |
| Central frequency | 150 MHz |
| Wavelength in ice | $\sim$1.12 m |
| Pulse duration (low gain) | 3 $\mu s$ |
| Pulse duration (high gain) | 10 $\mu s$ |
| Bandwidth | 10 MHz |
| Range resolution in ice (pulse compressed) | 8.4 m |
| Along-track resolution | 18 m, 30 m |



Fig. 7. Division of the training and validation datasets for our model.

ensuring randomness, unbiasedness, and representativeness in the sample distribution. Furthermore, the spatial distribution of the training and validation datasets was examined to ensure uniformity, allowing the validation dataset to effectively assess the model's performance across different regions in Gamburtsev area. As displayed in Fig. 7, the gray survey lines indicate the regions from which training samples are derived, while the red survey lines represent the validation dataset employed to evaluate model performance.

### A. Quantitative Analysis

To rigorously assess the performance of our proposed method, especially in the context of our highly imbalanced dataset, we utilize two metrics: $F1$ score and the distance

TABLE II
PERFORMANCE OF VARIOUS MODELS UNDER DIFFERENT POSITIVE-TO-NEGATIVE SAMPLE TRAINING RATIOS, I.E., 1:10 (POSITIVE: 1927 AND NEGATIVE: 192 70); 1:50 (POSITIVE: 1927 AND NEGATIVE: 96 350); AND 1:100 (POSITIVE: 1927 AND NEGATIVE: 192 700)

| | | ET | LightGBM | RF | CatBoost | Ours |
|---|---|---|---|---|---|---|
| F1 Score | 1:10 | 0.8185 | 0.7947 | 0.8152 | 0.7944 | 0.8195 |
| | 1:50 | 0.8016 | 0.7734 | 0.7983 | 0.7720 | 0.8192 |
| | | -1.69% | -2.13% | -1.69% | -2.23% | -0.03% |
| | 1:100 | 0.8000 | 0.7639 | 0.7959 | 0.7595 | 0.8163 |
| | | -1.85% | -3.08% | -1.93% | -3.49% | -0.32% |
| DIP | 1:10 | 0.8178 | 0.7944 | 0.8141 | 0.7940 | 0.8187 |
| | 1:50 | 0.8012 | 0.7727 | 0.7976 | 0.7711 | 0.8183 |
| | | -1.66% | -2.17% | -1.65% | -2.29% | -0.04% |
| | 1:100 | 0.7996 | 0.7626 | 0.7949 | 0.7583 | 0.8153 |
| | | -1.82% | -3.18% | -1.92% | -3.57% | -0.34% |

from the ideal position (DIP) [56], to evaluate the classification effectiveness of our method.

The aim of our study is to understand the impact of class imbalance on model accuracy. To address this, we created training datasets with three different positive-to-negative sample ratios: 1:10, 1:50, and 1:100, as listed in Table II. The value "1" (i.e., 1927) represents the total number of A-scope traces extracted from the gray survey lines in Fig. 7, corresponding to SL inventories identified by Livingstone et al. [13] and Wolovick et al. [18]. The negative samples, with counts equal to 10, 50, and 100 times the positive samples, were randomly selected from A-scope traces in regions of the gray survey lines where no SLs are present.

Table II presents the $F1$ score and DIP of base classifiers and our model under these different positive-to-negative sample ratios. The percentages indicate the decline rate in the metrics compared with the positive-to-negative sample ratio of 1:10.

The experimental results show that at a mild imbalanced sample ratio of 1:10, all models exhibited relatively high $F1$ scores. Among them, our model performed the best with an $F1$ score of 0.8195. However, as the sample imbalance ratio increased to 1:50 and 1:100, the performance of all models declined. Notably, the $F1$ score of our model decreased by only 0.03% at a 1:50 ratio and by 0.32% under the extreme imbalance condition of 1:100. Compared with other models, this performance decline is significantly smaller. In addition, the DIP value, another important indicator for assessing a model's ability to handle sample imbalances, also showed minimal reductions in our model. At a 1:10 sample ratio, our model's DIP value was 0.8187, with decreases of only 0.04% and 0.34% at 1:50 and 1:100 ratios, respectively. This further demonstrates the stability of our model in handling various imbalanced sample scenarios. These findings underscore the significant advantages of our model in countering sample imbalance in practical applications, compared with other models as a collective benchmark.

## B. Qualitative Analysis

In this section, we present the qualitative analysis of the automated prediction of SLs by our stacking model and showcase their distribution in the AGAP region.



Fig. 8. SL prediction results for various radar profiles within the AGAP region validation dataset. Red strips represent predicted SLs using our model, while blue strips represent SLs in the inventories. (a) L-1 and W-19. (b) W-39. (c) L-12. (d) L-37.

Fig. 8 presents four examples of prediction results for different radar profiles within our validation dataset in the AGAP region. For SLs predicted by Wolovick et al. [18], we use the notation "W-*," where "*" represents the corresponding lake's number in Wolovick's list. For lakes identified by Livingstone et al. [13], we use "L-*," where "*" indicates the lake's number in their list.

Specifically, Fig. 8(a)–(c) displays the comparison between our predicted results and the lakes previously identified by Livingstone et al. [13] and Wolovick et al. [18]. It is evident from the figures that the locations and lengths of the lakes predicted by our model closely match those listed in the existing records, visually confirming their accuracy. In Fig. 8(d),

Fig. 9. Examples of newly predicted SLs include (a) C-43, (b) C-1, and (c) C-16 and C-17.



Fig. 10. (a) Prediction result and smoothing process of L-27. (b) Slope value.



Fig. 11. Prediction result after smoothing process of L-31 and L-32.

lake L-37 is located in the northwestern part of the AGAP region, where the terrain appears relatively flat and the BRP is constantly higher than the adjacent regions. However, the location reported by Livingstone et al. [13] is slightly offset judged from the radargram, and from the length perspective, our result is more conservative and only the central region is predicted as the lake.

Next, we present the newly discovered candidate SLs identified by our model in regions outside of the original training and validation areas in the AGAP region. We also visualize the extracted features used as input for our stacking ensemble learning method, illustrating their contributions to the final prediction outcomes for the processed radargram frames. We denote these lakes as "C-*," where "*" corresponds to the lake's specific number in our list. The provided list of candidate SLs is evaluated based on CBRP values. The evaluation criterion is that if the difference in mean CBRP value between the current prediction range and any of the adjacent regions is greater than or equal to 15 dB [18], [57], [58], [59], the predicted lake can be listed. This additional step

Fig. 12. Prediction results for L-34, L-35, L-36, and newly predicted candidate SL C-50.



Fig. 13. SLs newly predicted in two adjacent B-scopes.

of screening is to ascertain the accuracy and reliability of our updated inventory, which is a conservative inventory.

Some examples of newly identified candidate SLs, which were not reported by Livingstone et al. [13] or Wolovick et al. [18], are presented in Fig. 9. The length of these candidate SLs, designated as C-43 [see Fig. 9(a)], C-1 [see Fig. 9(b)], C-16, and C-17 [see Fig. 9(c)], is approximately 558 m, 5 km, 1530 m, and 780 m, respectively. Based on the feature values and the traditional visual judgment of radargrams, all four lakes exhibit flat topography, high reflectivity, and elevated TFF values, with the CBRP exceeding that of the surrounding areas by more than 15 dB.

### C. Discussion

*1) Smoothing of the Prediction Results Based on A-Scope:* Since our method is based on A-scope level recognition, the predicted results may exhibit discontinuities, as shown in

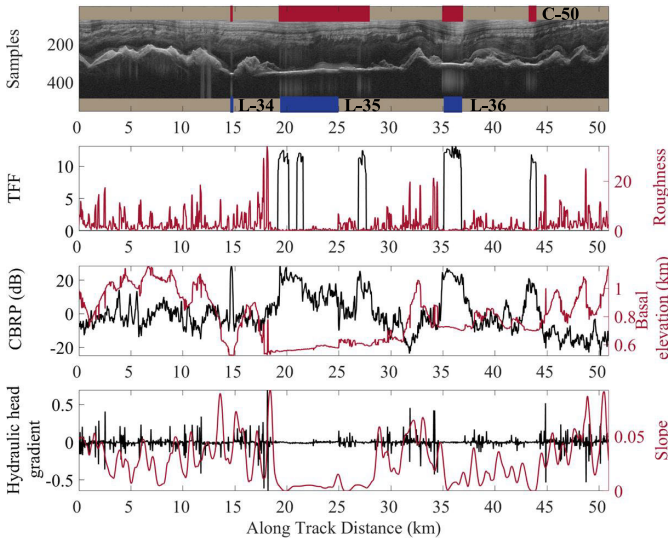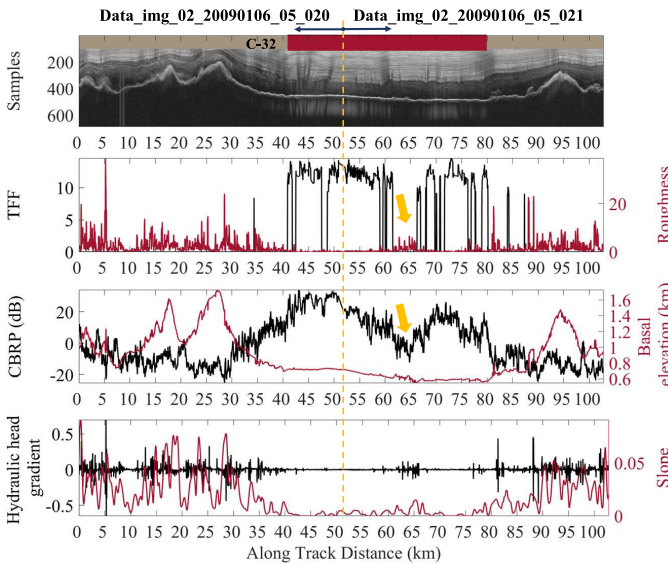Fig. 10(a). To address this issue, we applied a least squares fitting to the subglacial topography and used the first-order coefficient as the slope value to capture terrain variations. If the difference between the maximum and minimum slope values within a given interval is less than the empirically determined threshold, the region was considered suitable for smoothing for an SL. The slope values are shown in Fig. 10(b). In principle, this smoothing method transforms the prediction of subglacial water bodies at the A-scope level to the prediction of SLs, which effectively handles the discontinuity issue. However, such a smoothing operation is inevitably based on an empirical threshold, which serves as one of the limitations of the model. In the current study, we determine this threshold with which the predicted subglacial water bodies can be smoothed to match the inventoried SLs with an accuracy of 100%. From the perspective of continuous imaging, this smoothing method effectively enhances the model's performance in handling discontinuities.

However, the smoothing process may cause two nearby SLs to merge into a single SL. As depicted in Fig. 11, L-31 and L-32 are two separate SLs in the inventory provided by Livingstone et al. [13], but they are interpreted as a single SL after the smoothing process. This misclassification is likely due to the relatively flat terrain between the two lakes, which increases the probability that the model perceives them as a single, contiguous water body. For this case, we do not recognize this SL as a newly predicted SL.

*2) SLs Newly Predicted in the Validation Dataset:* In our predicted results, some newly predicted SLs are located on the survey line in the validation dataset. As shown in Fig. 12, SL C-50, predicted by our model, is not included in the current inventories. When accessing the accuracy of our prediction results, lake C-50 is classified as false positives, which may lower our $F1$ score and DIP.

*3) SLs Predicted in Two Adjacent B-Scopes:* In our results, some SLs are predicted at the ends of adjacent B-scopes, which likely represent the same lake. As shown in Fig. 13, we concatenate the two B-scope files, marking the boundary with a dashed line, resulting in a candidate SL of approximately 38 km in length. Overall, the lake exhibits pronounced CBRP characteristics and stable TFF. However, between 60 and 66 km, there is a decrease in CBRP, an increase in topographic roughness, and zero response in TFF (marked by yellow arrows). Nevertheless, our smoothing algorithm indicates that the terrain is relatively flat on a macroscopic scale. It may suggest the presence of a moist ice–bedrock interface, which might lead to speculation that this section may still be part of the same shallow lake. This interesting phenomenon suggests that future research should integrate additional data to conduct a thorough analysis of the geological and environmental characteristics of this area, to validate our current interpretations and to explore other possible geographical or climatic influences.

*4) Updated SL Inventory in the AGAP Region:* Wolovick et al. [18] produced an inventory for SLs in the AGAP region in 2013, which was included and updated in the worldwide SL inventory by Livingstone et al. [13]. The SLs in these two inventories are manually predicted, and

Fig. 14. Distribution of SLs in the updated inventory.



Fig. 15. Example of an undetected SL (L-42) in the validation dataset.



Fig. 16. Histogram of SL lengths in the updated inventory.

together with the re-examined SLs in [13] and [18], are listed in Table III for the AGAP region. In Table III, we show the coordinates of the center of the SL on the RES survey line and the length of the SL. Fig. 14 illustrates the spatial distribution of all SLs included in our updated SL inventory, as detailed in Table III. To enhance the clarity of regions where SLs appear densely clustered or overlapping in the main map, three inset panels—(I)–(III)—are provided to magnify these specific areas. The SLs previously reported by Livingstone et al. [13] and Wolovick et al. [18] are marked with yellow triangles and green pentagrams, respectively, while the newly identified candidate SLs in our study, which have not been reported before, are shown as blue circles. Furthermore, in Fig. 14, we use dashed circles to mark the locations of the only three lakes missed in the validation dataset, namely, L-26, L-42, and L-44. Taking L-42 as an example (as shown in Fig. 15), its CBRP contrast to neighboring areas is lower than that of the predicted lake C-52, and the local topographic slope is not prominent. These factors may collectively contribute to its omission by our algorithm.

it is found that there are certain overlaps between these inventories for the AGAP region. Through this study, our proposed method advances the traditional SL prediction by using an automated prediction model based on optimized stacking ensemble learning. Its capability in handling extreme data imbalance situations is validated for the AGAP region application, which has the lake-to-nonlake ratio approaching 1:100. We have predicted a total of 55 new candidate SLs with lengths ranging from 108 to 38 130 m. A new inventory is generated with a positive-to-negative sample ratio of 1:100. Noting that the actual positive-to-negative ratio in AGAP is 1:164, this reflects a much more realistic lake-to-nonlake ratio, compared with existing works [22], [31]. These candidate SLs,

TABLE III
INVENTORY OF SLs IN THE AGAP REGION OF ANTARCTICA

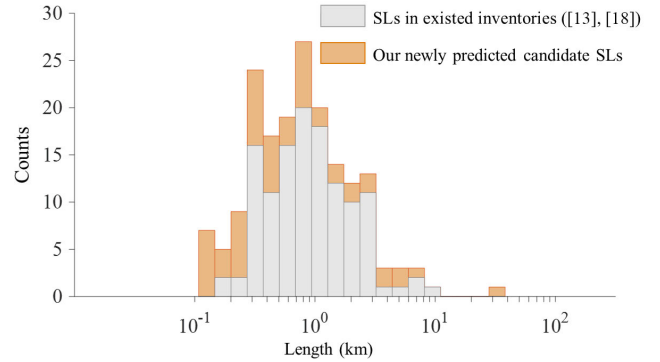| ID | Lake Name | File Name | Latitude | Longitude | Length (m) | Resolution (m) |
|---|---|---|---|---|---|---|
| 1 | L-1 | Data_img_02_20081223_01_006 | -83.2861 | 71.1630 | 1880 | 18 |
| 2 | L-2 | Data_img_02_20081223_01_007 | -82.9580 | 71.5740 | 7240 | 18 |
| 3 | L-3 | Data_img_02_20081223_01_007 | -82.8890 | 71.6550 | 990 | 18 |
| 4 | L-4 | Data_img_02_20081223_01_009 | -82.0390 | 72.5390 | 350 | 18 |
| 5 | L-5 | Data_img_02_20081223_01_014 | -79.9050 | 74.1080 | 180 | 18 |
| 6 | L-6 | Data_img_02_20081223_01_022 | -82.2260 | 74.0255 | 1330 | 18 |
| 7 | L-7 | Data_img_02_20081225_02_024 | -83.0400 | 69.9760 | 1490 | 30 |
| 8 | L-8 | Data_img_02_20081225_02_024 | -83.0630 | 69.9440 | 1020 | 30 |
| 9 | L-9 | Data_img_02_20081225_02_024 | -83.1010 | 69.8870 | 730 | 30 |
| 10 | L-10 | Data_img_02_20081225_02_024 | -83.1480 | 69.8200 | 580 | 30 |
| 11 | L-11 | Data_img_02_20081225_04_018 | -81.3270 | 77.9170 | 1370 | 30 |
| 12 | L-12 | Data_img_02_20090102_03_011 | -80.3300 | 80.0000 | 1070 | 18 |
| 13 | L-13 | Data_img_02_20090102_03_015 | -80.3020 | 81.0650 | 1210 | 18 |
| 14 | L-14 | Data_img_02_20090105_02_008 | -79.8980 | 71.2250 | 610 | 30 |
| 15 | L-15 | Data_img_02_20090106_05_009 | -80.9450 | 71.6260 | 300 | 30 |
| 16 | L-16 | Data_img_02_20090106_05_017 | -79.9660 | 59.4450 | 3080 | 30 |
| 17 | L-17 | Data_img_02_20090107_03_015 | -82.3960 | 72.7120 | 820 | 30 |
| 18 | L-18 | Data_img_02_20090107_03_015 | -82.4200 | 74.3590 | 1460 | 30 |
| 19 | L-19 | Data_img_02_20090107_03_015 | -82.4250 | 74.7131 | 2630 | 30 |
| 20 | L-20 | Data_img_02_20090107_05_020 | -80.9640 | 65.5070 | 1460 | 30 |
| 21 | L-21 | Data_img_02_20090108_01_016 | -80.2800 | 80.7960 | 2340 | 30 |
| 22 | L-22 | Data_img_02_20090108_01_016 | -80.4240 | 80.8080 | 880 | 30 |
| 23 | L-23 | Data_img_02_20090108_01_017 | -80.6780 | 80.8300 | 1170 | 30 |
| 24 | L-24 | Data_img_02_20090108_03_021 | -82.2630 | 74.3350 | 730 | 30 |
| 25 | L-25 | Data_img_02_20090108_03_021 | -82.3990 | 74.2340 | 780 | 30 |
| 26 | L-26 | Data_img_02_20090109_01_018 | -80.7910 | 74.1070 | 420 | 30 |
| 27 | L-27 | Data_img_02_20090109_04_005 | -83.2310 | 70.4630 | 2920 | 30 |
| 28 | L-28 | Data_img_02_20090109_04_024 | -83.2060 | 72.4140 | 2480 | 30 |
| 29 | L-29 | Data_img_02_20090110_01_008 | -82.1180 | 81.9590 | 1800 | 30 |
| 30 | L-30 | Data_img_02_20090110_01_008 | -81.9470 | 81.9180 | 1520 | 30 |
| 31 | L-31 | Data_img_02_20090110_01_009 | -81.5840 | 81.8370 | 180 | 30 |
| 32 | L-32 | Data_img_02_20090110_01_009 | -81.5740 | 81.8350 | 670 | 30 |
| 33 | L-33 | Data_img_02_20090111_02_025 | -82.4070 | 72.0420 | 270 | 18 |
| 34 | L-34 | Data_img_02_20090106_05_013 | -80.3710 | 66.2579 | 290 | 30 |
| 35 | L-35 | Data_img_02_20090106_05_013 | -80.3552 | 65.8700 | 5610 | 30 |
| 36 | L-36 | Data_img_02_20090106_05_013 | -80.3248 | 65.1567 | 1760 | 30 |
| 37 | L-37 | Data_img_02_20090109_04_004 | -83.3956 | 70.2276 | 8310 | 30 |
| 38 | L-39 | Data_img_02_20090108_03_005 | -82.9679 | 71.9315 | 1750 | 30 |
| 39 | L-40 | Data_img_02_20090108_03_021 | -82.4606 | 74.1864 | 730 | 30 |
| 40 | L-41 | Data_img_02_20081228_01_006 | -82.5455 | 75.1611 | 1980 | 30 |
| 41 | L-42 | Data_img_02_20090108_01_007 | -82.2566 | 80.3314 | 480 | 30 |
| 42 | L-43 | Data_img_02_20090108_01_009 | -81.7247 | 80.3187 | 1020 | 30 |
| 43 | L-44 | Data_img_02_20090102_03_015 | -80.3345 | 81.0697 | 1005 | 18 |
| 44 | W-1 | Data_img_02_20081225_02_014 | -79.7050 | 71.9570 | 740 | 30 |
| 45 | W-2 | Data_img_02_20081229_01_026 | -83.9400 | 66.7540 | 1780 | 18 |
| 46 | W-3 | Data_img_02_20081229_01_018 | -80.7890 | 71.2840 | 1170 | 18 |
| 47 | W-4 | Data_img_02_20081229_01_007 | -82.8890 | 69.4530 | 1250 | 18 |
| 48 | W-5 | Data_img_02_20081229_01_010 | -81.6730 | 70.9900 | 640 | 18 |
| 49 | W-6 | F45a_L320-175_HGe4 | -83.1550 | 69.4240 | 610 | 30 |
| 50 | W-7 | F45a_L320-182_HGe4 | -82.9490 | 69.7350 | 3050 | 30 |
| 51 | W-8 | F45a_L320-181_HGe4 | -82.9730 | 69.7000 | 630 | 30 |
| 52 | W-12 | Data_img_02_20081228_03_026 | -83.1800 | 70.1540 | 2570 | 18 |
| 53 | W-13 | Data_img_02_20081228_03_024 | -82.3180 | 71.2510 | 330 | 18 |

TABLE III
(*Continued.*) INVENTORY OF SLs IN THE AGAP REGION OF ANTARCTICA

| ID | Lake Name | File Name | Latitude | Longitude | Length (m) | Resolution (m) |
|---|---|---|---|---|---|---|
| 54 | W-16 | Data_img_02_20081228_03_006 | -83.3100 | 70.7460 | 7500 | 18 |
| 55 | W-17 | Data_img_02_20081228_03_006 | -83.2590 | 70.8180 | 710 | 18 |
| 56 | W-18 | Data_img_02_20081228_03_007 | -82.8820 | 71.2960 | 2710 | 18 |
| 57 | W-19 | Data_img_02_20081223_01_006 & 007 | -82.9910 | 71.5360 | 1210 | 18 |
| 58 | W-21 | Data_img_02_20081223_01_009 & 010 | -81.7150 | 72.8280 | 230 | 18 |
| 59 | W-22 | Data_img_02_20081223_01_012 | -80.8450 | 73.5050 | 350 | 18 |
| 60 | W-23 | Data_img_02_20081223_01_014 | -79.9040 | 74.1080 | 350 | 18 |
| 61 | W-25 | Data_img_02_20090109_01_023 | -83.0000 | 72.2630 | 280 | 30 |
| 62 | W-26 | Data_img_02_20090109_01_018 | -80.8610 | 74.0630 | 1940 | 30 |
| 63 | W-28 | Data_img_02_20090109_01_012 | -79.9480 | 75.1140 | 390 | 30 |
| 64 | W-30 | Data_img_02_20081223_01_021 | -82.1830 | 74.0630 | 1910 | 18 |
| 65 | W-34 | Data_img_02_20081228_01_020 | -82.3960 | 74.5760 | 2640 | 30 |
| 66 | W-35 | Data_img_02_20081228_01_020 | -82.3390 | 74.6150 | 1170 | 30 |
| 67 | W-36 | Data_img_02_20081227_01_023 | -82.9970 | 74.4840 | 320 | 30 |
| 68 | W-37 | Data_img_02_20081227_01_022 | -82.4650 | 74.8690 | 850 | 30 |
| 69 | W-38 | Data_img_02_20081227_01_022 | -82.4500 | 74.8790 | 670 | 30 |
| 70 | W-39 | Data_img_02_20081228_01_009 | -81.0020 | 75.9860 | 430 | 30 |
| 71 | W-40 | Data_img_02_20081231_01_007 | -82.2880 | 76.9940 | 1450 | 30 |
| 72 | W-41 | Data_img_02_20081231_01_008 | -81.7950 | 77.1740 | 640 | 30 |
| 73 | W-42 | Data_img_02_20081225_04_020 | -82.3470 | 77.6440 | 510 | 30 |
| 74 | W-43 | F51c_L540-309_HGe4 | -82.4250 | 77.9630 | 2410 | 30 |
| 75 | W-44 | F51c_L540-310_HGe4 | -82.4120 | 77.9660 | 1330 | 30 |
| 76 | W-45 | F51c_L540-308_HGe4 | -82.3980 | 76.9690 | 1780 | 30 |
| 77 | W-46 | Data_img_02_20090113_01_005 | -82.5490 | 78.2730 | 3640 | 30 |
| 78 | W-47 | Data_img_02_20090113_01_006 | -82.5070 | 78.2820 | 980 | 30 |
| 79 | W-48 | Data_img_02_20090113_01_006 | -82.4920 | 78.2850 | 870 | 30 |
| 80 | W-49 | Data_img_02_20090113_01_006 | -82.4530 | 78.2950 | 1070 | 30 |
| 81 | W-50 | Data_img_02_20090113_01_008 | -81.2580 | 78.5260 | 370 | 30 |
| 82 | W-51 | F15b_L580-252_HGe4 | -81.9180 | 79.3620 | 440 | 30 |
| 83 | W-52 | Data_img_02_20090110_03_023 | -83.7030 | 79.5950 | 450 | 30 |
| 84 | W-53 | Data_img_02_20090102_03_008 | -81.8180 | 79.9990 | 310 | 18 |
| 85 | W-55 | Data_img_02_20090108_01_008 | -81.7460 | 80.3120 | 790 | 30 |
| 86 | W-57 | F42a_L620-085_HGe4 | -82.8880 | 80.7210 | 1090 | 30 |
| 87 | W-58 | Data_img_02_20090108_01_019 | -81.7400 | 80.9350 | 790 | 30 |
| 88 | W-61 | Data_img_02_20090102_03_023 | -83.6740 | 81.6230 | 610 | 18 |
| 89 | W-62 | Data_img_02_20090102_03_021 | -82.7950 | 81.4270 | 570 | 18 |
| 90 | W-63 | Data_img_02_20090102_03_016 | -80.6280 | 81.1000 | 310 | 18 |
| 91 | W-65 | F18a_L650-211_HGe4 | -82.1500 | 81.6410 | 400 | 30 |
| 92 | W-68 | Data_img_02_20090110_01_008 | -81.9500 | 81.9190 | 920 | 30 |
| 93 | W-69 | Data_img_02_20090110_01_008 | -82.0060 | 81.9330 | 670 | 30 |
| 94 | W-73 | Data_img_02_20090105_01_007 | -82.6970 | 82.4700 | 1420 | 30 |
| 95 | W-74 | Data_img_02_20090105_01_010 | -81.5410 | 82.1360 | 540 | 30 |
| 96 | W-75 | Data_img_02_20090105_01_014 | -79.8770 | 81.7860 | 1250 | 30 |
| 97 | W-76 | F42b_L690-262_HGe4 | -81.4140 | 82.7030 | 920 | 30 |
| 98 | W-77 | Data_img_02_20090103_02_020 | -81.3780 | 82.9930 | 430 | 18 |
| 99 | W-78 | Data_img_02_20090103_02_010 | -81.5120 | 84.5600 | 400 | 18 |
| 100 | W-79 | Data_img_02_20090103_02_012 | -80.9110 | 84.2590 | 910 | 18 |
| 101 | W-80 | Data_img_02_20090105_02_019 | -81.5160 | 84.8670 | 980 | 30 |
| 102 | W-81 | F52a_T10120-126_HGe4 | -83.2480 | 69.3270 | 1010 | 30 |
| 103 | W-82 | F52a_T10120-120_HGe4 | -83.2770 | 70.7460 | 1630 | 30 |
| 104 | W-83 | F52a_T10120-119_HGe4 | -83.2820 | 71.0390 | 2390 | 30 |
| 105 | W-84 | F52b_T10130-153_HGe4 | -82.9490 | 69.7030 | 1240 | 30 |
| 106 | W-85 | F52b_T10130-160_HGe4 | -82.9850 | 71.4680 | 520 | 30 |
| 107 | W-86 | F52b_T10130-162_HGe4 | -82.9980 | 72.2330 | 910 | 30 |
| 108 | W-87 | F52b_T10130-194_HGe4 | -83.0600 | 80.9720 | 2380 | 30 |

TABLE III
(*Continued.*) Inventory of SLs in the AGAP Region of Antarctica

| ID | Lake Name | File Name | Latitude | Longitude | Length (m) | Resolution (m) |
|---|---|---|---|---|---|---|
| 109 | W-91 | Data_img_02_20090107_03_015 & 016 | -82.4260 | 74.7750 | 320 | 30 |
| 110 | W-92 | Data_img_02_20090107_03_016 & 017 | -82.4530 | 78.2360 | 720 | 30 |
| 111 | W-93 | Data_img_02_20090107_03_017 | -82.4530 | 78.2810 | 340 | 30 |
| 112 | W-94 | Data_img_02_20090107_03_016 & 017 | -82.4530 | 78.1640 | 500 | 30 |
| 113 | W-95 | Data_img_02_20090107_03_016 | -82.4530 | 78.1320 | 490 | 30 |
| 114 | W-96 | Data_img_02_20090107_03_016 | -82.4520 | 78.1070 | 290 | 30 |
| 115 | W-97 | Data_img_02_20090107_03_016 | -82.4520 | 78.0740 | 650 | 30 |
| 116 | W-98 | Data_img_02_20090107_03_016 | -82.4520 | 78.0380 | 290 | 30 |
| 117 | W-99 | Data_img_02_20090106_04_010 | -82.1520 | 81.5590 | 890 | 30 |
| 118 | W-100 | Data_img_02_20090106_04_010 | -82.1510 | 81.6330 | 310 | 30 |
| 119 | W-101 | Data_img_02_20090106_04_010 | -82.1500 | 81.7830 | 930 | 30 |
| 120 | W-102 | Data_img_02_20090111_02_013 | -81.5430 | 82.1220 | 550 | 18 |
| 121 | W-103 | Data_img_02_20090111_02_014 | -81.5190 | 84.8920 | 1500 | 18 |
| 122 | W-104 | Data_img_02_20090109_03_015 | -80.9190 | 84.2730 | 1550 | 30 |
| 123 | W-105 | Data_img_02_20090109_03_018 | -80.6390 | 81.0220 | 790 | 30 |
| 124 | C-1 | Data_img_02_20081223_01_015 | -79.5831 | 74.5711 | 4932 | 18 |
| 125 | C-2 | Data_img_02_20081225_02_028 | -84.3856 | 73.4685 | 360 | 30 |
| 126 | C-3 | Data_img_02_20081226_02_004 | -83.6404 | 82.4832 | 1998 | 18 |
| 127 | C-4 | Data_img_02_20081226_02_026 | -83.2581 | 82.5699 | 990 | 18 |
| 128 | C-5 | Data_img_02_20081227_01_010 | -80.7354 | 76.3785 | 120 | 30 |
| 129 | C-6 | Data_img_02_20081228_03_011 | -81.3163 | 72.8579 | 108 | 18 |
| 130 | C-7 | Data_img_02_20081229_01_009 | -82.1214 | 70.4773 | 108 | 18 |
| 131 | C-8 | Data_img_02_20081229_01_016 | -79.7886 | 72.1339 | 2484 | 18 |
| 132 | C-9 | Data_img_02_20081231_01_005 | -83.0351 | 76.6761 | 360 | 30 |
| 133 | C-10 | Data_img_02_20081231_01_009 | -81.2581 | 77.3368 | 390 | 30 |
| 134 | C-11 | Data_img_02_20081231_01_010 | -80.8103 | 77.4741 | 120 | 30 |
| 135 | C-12 | Data_img_02_20081231_01_018 | -80.7524 | 76.9207 | 270 | 30 |
| 136 | C-13 | Data_img_02_20081231_01_019 | -80.6125 | 76.6185 | 630 | 30 |
| 137 | C-14 | Data_img_02_20081231_03_006 | -83.1558 | 69.3630 | 450 | 30 |
| 138 | C-15 | Data_img_02_20081231_03_009 | -82.8176 | 64.9017 | 780 | 30 |
| 139 | C-16 | Data_img_02_20081231_03_011 | -83.2504 | 65.4084 | 1530 | 30 |
| 140 | C-17 | Data_img_02_20081231_03_011 | -83.2986 | 65.7526 | 780 | 30 |
| 141 | C-18 | Data_img_02_20090101_01_008 | -81.8189 | 80.0071 | 144 | 18 |
| 142 | C-19 | Data_img_02_20090101_01_008 | -81.7088 | 79.9994 | 108 | 18 |
| 143 | C-20 | Data_img_02_20090103_03_018 | -80.8617 | 83.6714 | 270 | 30 |
| 144 | C-21 | Data_img_02_20090103_03_019 | -81.0642 | 83.7503 | 360 | 30 |
| 145 | C-22 | Data_img_02_20090103_03_019 | -81.2010 | 83.8107 | 240 | 30 |
| 146 | C-23 | Data_img_02_20090105_01_005 | -83.2787 | 82.6752 | 690 | 30 |
| 147 | C-24 | Data_img_02_20090105_02_017 | -80.9324 | 84.5547 | 1770 | 30 |
| 148 | C-25 | Data_img_02_20090105_02_018 | -81.0727 | 84.6264 | 150 | 30 |
| 149 | C-26 | Data_img_02_20090105_03_005 | -83.7178 | 79.3174 | 5040 | 30 |
| 150 | C-27 | Data_img_02_20090105_03_012 | -80.8816 | 83.8409 | 240 | 30 |
| 151 | C-28 | Data_img_02_20090106_01_012 | -81.2236 | 83.5261 | 300 | 30 |
| 152 | C-29 | Data_img_02_20090106_04_004 | -83.2407 | 82.4197 | 690 | 30 |
| 153 | C-30 | Data_img_02_20090106_04_026 | -83.7919 | 80.1360 | 1020 | 30 |
| 154 | C-31 | Data_img_02_20090106_05_019 | -79.5730 | 62.6904 | 4260 | 30 |
| 155 | C-32 | Data_img_02_20090106_05_020 & 21 | -79.7214 | 65.6525 | 38130 | 30 |
| 156 | C-33 | Data_img_02_20090107_03_002 | -83.7113 | 81.7111 | 420 | 30 |
| 157 | C-34 | Data_img_02_20090107_05_010 | -81.4677 | 62.5884 | 1440 | 30 |
| 158 | C-35 | Data_img_02_20090108_03_002 | -84.3429 | 72.8579 | 150 | 30 |
| 159 | C-36 | Data_img_02_20090108_03_007 | -82.4095 | 72.5205 | 480 | 30 |
| 160 | C-37 | Data_img_02_20090109_03_005 | -82.4465 | 72.5499 | 120 | 30 |
| 161 | C-38 | Data_img_02_20090109_03_026 | -82.5863 | 75.4818 | 3540 | 30 |
| 162 | C-39 | Data_img_02_20090110_01_011 | -81.0189 | 81.7243 | 360 | 30 |
| 163 | C-40 | Data_img_02_20090110_01_018 | -81.3963 | 82.3929 | 300 | 30 |

TABLE III
*(Continued.)* INVENTORY OF SLs IN THE AGAP REGION OF ANTARCTICA

| ID | Lake Name | File Name | Latitude | Longitude | Length (m) | Resolution (m) |
|-----|-----------|--------------------------------|----------|-----------|------------|----------------|
| 164 | C-41 | Data_img_02_20090110_01_018 | -81.4440 | 82.4085 | 180 | 30 |
| 165 | C-42 | Data_img_02_20090110_03_010 | -80.6137 | 79.1766 | 510 | 30 |
| 166 | C-43 | Data_img_02_20090111_02_006 | -82.3189 | 71.2807 | 558 | 18 |
| 167 | C-44 | Data_img_02_20090111_02_007 | -81.8197 | 70.3866 | 342 | 18 |
| 168 | C-45 | Data_img_02_20090111_02_011 | -81.5372 | 76.9461 | 324 | 18 |
| 169 | C-46 | Data_img_02_20090111_02_016 | -81.2315 | 83.3768 | 846 | 18 |
| 170 | C-47 | Data_img_02_20081228_03_026 | -83.1099 | 70.2539 | 234 | 18 |
| 171 | C-48 | Data_img_02_20081229_01_007 | -82.8431 | 69.5211 | 864 | 18 |
| 172 | C-49 | Data_img_02_20090103_02_020 | -81.4288 | 83.0068 | 252 | 18 |
| 173 | C-50 | Data_img_02_20090106_05_013 | -80.3074 | 64.7632 | 750 | 30 |
| 174 | C-51 | Data_img_02_20090106_05_017 | -79.8636 | 59.6943 | 7950 | 30 |
| 175 | C-52 | Data_img_02_20090108_01_007 | -82.3138 | 80.3341 | 210 | 30 |
| 176 | C-53 | Data_img_02_20090108_01_008 | -81.8300 | 80.3145 | 480 | 30 |
| 177 | C-54 | Data_img_02_20090109_01_023 | -83.1518 | 72.0934 | 2550 | 30 |
| 178 | C-55 | Data_img_02_20090113_01_008 | -81.3475 | 78.5091 | 420 | 30 |

Fig. 16 shows the length distribution of newly predicted SLs and those in the existing inventories [13], [18]. In the AGAP region, the lengths of SLs are primarily concentrated between 300 and 3000 m, with most newly predicted SLs measuring less than 1 km. This phenomenon is closely related to the topographical features of the AGAP region. The rugged terrain of the Gamburtsev Mountains facilitates the accumulation of water in small bedrock depressions, leading to a prevalence of small lakes in these low-lying areas [18]. On the other hand, a rather large lake is predicted, as displayed in Fig. 13. It measures 38 130 m in length, although we speculate that this SL may be segmented by wet sediment in the middle.

## VI. CONCLUSION AND FUTURE WORK

In this study, we proposed an automated method for predicting SLs. The novelty of this approach lies in addressing the inherent imbalance issue in RES data using an ABC-assisted stacking ensemble learning method, which aids in selecting the optimum classifier combination. We implemented the proposed method in the AGAP region of East Antarctica, conducting a comprehensive analysis that includes both quantitative and qualitative assessments. It is found that the performance of the proposed SLs prediction method is much more stable than the traditional models, as the imbalance between positive and negative samples increases. Another advantage of our method is its prediction capability at the A-scope level with the highest possible along-track resolution, demonstrating its adaptability and precision across various SL sizes, including the ability to predict small SLs. Finally, we generated an updated SL inventory for the AGAP region, predicting a total of 55 new candidate SLs with lengths ranging from 108 to 38 130 m. The spatial and temporal distribution of these new SLs and their potential hydraulic connections present an intriguing direction for future research. We also plan to further investigate the impact of basal roughness on radar reflection signals and explore how significant changes in RES data are commonly associated with

changes in bedrock wetness [44], [45]. It is worth investigating the temperature-dependent variations in the permittivity of the same material, as well as the relationship between the dielectric permittivity of materials and their moisture content. In addition, we will incorporate seismology or other geophysical data to more accurately assess the existence of SLs [60], [61].

## REFERENCES

[1] S. Willcocks and D. Hasterok, "Prediction of subglacial lake melt source regions from site characteristics," *Antarctic Sci.*, vol. 35, no. 2, pp. 127–140, Apr. 2023.

[2] E. MacKie, D. Schroeder, J. Caers, M. Siegfried, and C. Scheidt, "Antarctic topographic realizations and geostatistical modeling used to map subglacial lakes," *J. Geophys. Res., Earth Surf.*, vol. 125, no. 3, 2020, Art. no. e2019JF005420.

[3] J. K. Ridley, W. Cudlip, and S. W. Laxon, "Identification of subglacial lakes using ERS-1 radar altimeter," *J. Glaciol.*, vol. 39, no. 133, pp. 625–634, 1993.

[4] M. J. Siegert, N. Ross, and A. M. Le Brocq, "Recent advances in understanding Antarctic subglacial lakes and hydrology," *Phil. Trans. Roy. Soc. A, Math., Phys. Eng. Sci.*, vol. 374, no. 2059, Jan. 2016, Art. no. 20140306.

[5] B. E. Smith, H. A. Fricker, I. R. Joughin, and S. Tulaczyk, "An inventory of active subglacial lakes in Antarctica detected by ICESat (2003–2008)," *J. Glaciol.*, vol. 55, no. 192, pp. 573–595, 2009.

[6] H. Björnsson, "Subglacial lakes and jökulhlaups in Iceland," *Global Planet. Change*, vol. 35, nos. 3–4, pp. 255–271, 2003.

[7] R. Dziadek, F. Ferraccioli, and K. Gohl, "High geothermal heat flow beneath Thwaites glacier in west Antarctica inferred from aeromagnetic data," *Commun. Earth Environ.*, vol. 2, no. 1, p. 162, Aug. 2021.

[8] H. Björnsson, "Jökulhlaups in Iceland: Prediction, characteristics and simulation," *Ann. Glaciol.*, vol. 16, pp. 95–106, Jan. 1992.

[9] A. Burton-Johnson, R. Dziadek, and C. Martin, "Review article: Geothermal heat flow in antarctica: Current and future directions," *Cryosphere*, vol. 14, no. 11, pp. 3843–3873, Nov. 2020.

[10] P. Friedl, F. Weiser, A. Fluhrer, and M. H. Braun, "Remote sensing of glacier and ice sheet grounding lines: A review," *Earth-Sci. Rev.*, vol. 201, Feb. 2020, Art. no. 102948.

[11] D. M. Schroeder, "Paths forward in radioglaciology," *Ann. Glaciol.*, vol. 63, nos. 87–89, pp. 13–17, Sep. 2022.

[12] D. M. Schroeder et al., "Five decades of radioglaciology," *Ann. Glaciol.*, vol. 61, no. 81, pp. 1–13, Apr. 2020.

[13] S. J. Livingstone et al., "Subglacial lakes and their changing role in a warming climate," *Nature Rev. Earth Environ.*, vol. 3, no. 2, pp. 106–124, Jan. 2022.

[14] G. Oswald and G. d. Q. Robin, "Lakes beneath the Antarctic ice sheet," *Nature*, vol. 245, no. 5423, pp. 251–254, 1973.

[15] M. Siegert, J. Dowdeswell, M. Gorman, and N. McIntyre, "An inventory of Antarctic sub-glacial lakes," *Antarctic Sci.*, vol. 8, no. 3, pp. 281–286, 1996.

[16] M. J. Siegert, S. P. Carter, I. E. Tabacco, S. V. Popov, and D. D. Blankenship, "A revised inventory of Antarctic subglacial lakes," *Antarctic Sci.*, vol. 17, no. 3, pp. 453–460, 2004.

[17] A. Wright and M. Siegert, "A fourth inventory of Antarctic subglacial lakes," *Antarctic Sci.*, vol. 24, no. 6, pp. 659–664, Dec. 2012.

[18] M. J. Wolovick, R. E. Bell, T. T. Creyts, and N. Frearson, "Identification and control of subglacial water networks under Dome A, Antarctica," *J. Geophys. Res., Earth Surf.*, vol. 118, no. 1, pp. 140–154, 2013.

[19] S. V. Popov and V. N. Masolov, "Forty-seven new subglacial lakes in the 0–110° e sector of east Antarctica," *J. Glaciology*, vol. 53, no. 181, pp. 289–297, 2007.

[20] M. J. Siegert, "Antarctic subglacial lakes," *Earth-Sci. Rev.*, vol. 50, nos. 1–2, pp. 29–50, May 2000.

[21] S. P. Carter, D. D. Blankenship, M. E. Peters, D. A. Young, J. W. Holt, and D. L. Morse, "Radar-based subglacial lake classification in Antarctica," *Geochem., Geophys., Geosyst.*, vol. 8, no. 3, pp. 1–20, 2007.

[22] S. Lang et al., "A semiautomatic method for predicting subglacial dry and wet zones through identifying Dry–Wet transitions," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, 2022, Art. no. 4306815.

[23] X. Li et al., "Characterization of subglacial landscapes by a two-parameter roughness index," *J. Glaciol.*, vol. 56, no. 199, pp. 831–836, 2010.

[24] G. Oswald and S. Gogineni, "Recovery of subglacial water extent from Greenland radar survey data," *J. Glaciol.*, vol. 54, no. 184, pp. 94–106, 2008.

[25] G. K. A. Oswald, S. Rezvanbehbahani, and L. A. Stearns, "Radar evidence of ponded subglacial water in Greenland," *J. Glaciol.*, vol. 64, no. 247, pp. 711–729, 2018.

[26] T. M. Jordan et al., "Self-affine subglacial roughness: Consequences for radar scattering and basal water discrimination in northern Greenland," *Cryosphere*, vol. 11, no. 3, pp. 1247–1264, May 2017.

[27] A. Zirizzotti, L. Cafarella, and S. Urbini, "Ice and bedrock characteristics underneath dome c (Antarctica) from radio echo sounding data analysis," *IEEE Trans. Geosci. Remote Sens.*, vol. 50, no. 1, pp. 37–43, Jan. 2012.

[28] M. A. Cooper, T. M. Jordan, D. M. Schroeder, M. J. Siegert, C. N. Williams, and J. L. Bamber, "Subglacial roughness of the Greenland ice sheet: Relationship with contemporary ice velocity and geology," *Cryosphere*, vol. 13, no. 11, pp. 3093–3115, Nov. 2019.

[29] T. Hao et al., "Automatic detection of subglacial water bodies in the AGAP region, east antarctica, based on short-time Fourier transform," *Remote Sens.*, vol. 15, no. 2, p. 363, Jan. 2023.

[30] S. Qian and D. Chen, "Joint time-frequency analysis," *IEEE Signal Process. Mag.*, vol. 16, no. 2, pp. 52–67, Mar. 1999.

[31] A.-M. Ilisei, M. Khodadadzadeh, A. Ferro, and L. Bruzzone, "An automatic method for subglacial lake detection in ice sheet radar sounder data," *IEEE Trans. Geosci. Remote Sens.*, vol. 57, no. 6, pp. 3252–3270, Jun. 2019.

[32] S. Dong, L. Fu, X. Tang, Z. Li, and X. Chen, "Deep clustering in radar subglacial reflector reveals new subglacial lakes," *Cryosphere Discuss.*, vol. 2023, pp. 1–24, Jun. 2023.

[33] V. Nasteski, "An overview of the supervised machine learning methods," *Horizons. B*, vol. 4, nos. 51–62, p. 56, 2017.

[34] Y. Wang, D. Wang, N. Geng, Y. Wang, Y. Yin, and Y. Jin, "Stacking-based ensemble learning of decision trees for interpretable prostate cancer detection," *Appl. Soft Comput.*, vol. 77, pp. 188–204, Apr. 2019.

[35] F. Haghighi and H. Omranpour, "Stacking ensemble model of deep learning and its application to Persian/Arabic handwritten digits recognition," *Knowl.-Based Syst.*, vol. 220, May 2021, Art. no. 106940.

[36] Z. Wang, M. Zheng, and P. X. Liu, "A novel classification method based on stacking ensemble for imbalanced problems," *IEEE Trans. Instrum. Meas.*, vol. 72, pp. 1–13, 2023.

[37] R. Shreve, "Movement of water in glaciers," *J. Glaciol.*, vol. 11, no. 62, pp. 205–214, 1972.

[38] M. J. Siegert, J. Taylor, and A. J. Payne, "Spectral roughness of subglacial topography and implications for former ice-sheet dynamics in East Antarctica," *Global Planet. Change*, vol. 45, nos. 1–3, pp. 249–263, 2005.

[39] S. Fujita, T. Matsuoka, T. Ishida, K. Matsuoka, and S. Mae, "A summary of the complex dielectric permittivity of ice in the megahertz range and its applications for radar sounding of polar ice sheets," in *Physics of Ice Core Records*. Hokkaido, Japan: Univ. Press, 2000, pp. 185–212.

[40] S. M. Tulaczyk and N. T. Foley, "The role of electrical conductivity in radar wave reflection from glacier beds," *Cryosphere*, vol. 14, no. 12, pp. 4495–4506, Dec. 2020.

[41] D. M. Schroeder, D. D. Blankenship, R. K. Raney, and C. Grima, "Estimating subglacial water geometry using radar bed echo specularity: Application to Thwaites glacier, west Antarctica," *IEEE Geosci. Remote Sens. Lett.*, vol. 12, no. 3, pp. 443–447, Mar. 2015.

[42] A. Rutishauser et al., "Radar sounding survey over Devon ice cap indicates the potential for a diverse hypersaline subglacial hydrological environment," *Cryosphere*, vol. 16, no. 2, pp. 379–395, Feb. 2022.

[43] Y. Li, Y. Lu, and M. J. Siegert, "Radar sounding confirms a hydrologically active deep-water subglacial lake in east Antarctica," *Frontiers Earth Sci.*, vol. 8, p. 294, Jul. 2020.

[44] A. L. Broome and D. M. Schroeder, "A radiometrically precise multi-frequency ice-penetrating radar architecture," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, 2022, Art. no. 5104515.

[45] K. C. Jezek, M. Brogioni, J. T. Johnson, D. M. Schroeder, A. L. Broome, and G. Macelloni, "Active and passive microwave remote sensing of priestley glacier, Antarctica," *IEEE Trans. Geosci. Remote Sens.*, vol. 62, 2024, Art. no. 4302708.

[46] S. Cui, Y. Yin, D. Wang, Z. Li, and Y. Wang, "A stacking-based ensemble learning method for earthquake casualty prediction," *Appl. Soft Comput.*, vol. 101, Mar. 2021, Art. no. 107038.

[47] N. Kardani, A. Zhou, M. Nazem, and S.-L. Shen, "Improved prediction of slope stability using a hybrid stacking ensemble method based on finite element analysis and field data," *J. Rock Mech. Geotechnical Eng.*, vol. 13, no. 1, pp. 188–201, Feb. 2021.

[48] D. Karaboga, "An idea based on honey bee swarm for numerical optimization," Eng. Fac., Comput. Eng. Dept., Erciyes Univ., Kayseri, Türkiye, Tech. Rep. tr06, 2005.

[49] T. Chen and C. Guestrin, "XGBoost: A scalable tree boosting system," in *Proc. 22nd ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining*, Aug. 2016, pp. 785–794.

[50] G. Ke et al., "LightGBM: A highly efficient gradient boosting decision tree," in *Proc. Adv. Neural Inf. Process. Syst.*, vol. 30, 2017, pp. 1–9.

[51] T. M. Khoshgoftaar, M. Golawala, and J. V. Hulse, "An empirical study of learning from imbalanced data using random forest," in *Proc. 19th IEEE Int. Conf. Tools Artif. Intell. (ICTAI)*, vol. 2, Oct. 2007, pp. 310–317.

[52] R. E. Bell et al., "Widespread persistent thickening of the east Antarctic ice sheet by freezing from the base," *Science*, vol. 331, no. 6024, pp. 1592–1595, Mar. 2011.

[53] H. Corr, F. Ferraccioli, T. Jordan, and C. Robinson, "Processed airborne radio-echo sounding data from the agap survey covering Antarctica's Gamburtsev Province, East Antarctica (2007/2009) (version 1.0) [data set]," NERC EDS UK Polar Data Centre, Cambridge, U.K., 2021. [Online]. Available: https://doi.org/10.5285/A1ABF071-85FC-4118-AD37-7F186B72C847 and https://data.bas.ac.uk/full-record.php?id=GB/NERC/BAS/PDC/01544

[54] CReSIS. (2016). *CReSIS Radar Depth Sounder Data*. Lawrence, Kansas, USA. Digital Media. [Online]. Available: http://data.cresis.ku.edu/

[55] F. Ferraccioli, C. A. Finn, T. A. Jordan, R. E. Bell, L. M. Anderson, and D. Damaske, "East Antarctic rifting triggers uplift of the Gamburtsev mountains," *Nature*, vol. 479, no. 7373, pp. 388–392, Nov. 2011.

[56] A. K. Nandi, "From multiple independent metrics to single performance measure based on objective function," *IEEE Access*, vol. 11, pp. 3899–3913, 2023.

[57] C. Pierce et al., "Characterizing sub-glacial hydrology using radar simulations," *Cryosphere*, vol. 18, no. 4, pp. 1495–1515, Apr. 2024. [Online]. Available: https://tc.copernicus.org/articles/18/1495/2024/

[58] A. Rutishauser et al., "Discovery of a hypersaline subglacial lake complex beneath Devon ice cap, Canadian Arctic," *Sci. Adv.*, vol. 4, no. 4, Apr. 2018, Art. no. eaar4353.

[59] K. Matsuoka, "Pitfalls in radar diagnosis of ice-sheet bed conditions: Lessons from englacial attenuation models," *Geophys. Res. Lett.*, vol. 38, no. 5, 2011, Art. no. L05505.

[60] C. Hofstede et al., "The subglacial lake that Wasn't there: Improved interpretation from seismic data reveals a sediment bedform at isunnguata sermia," *J. Geophys. Res., Earth Surf.*, vol. 128, no. 10, Oct. 2023, Art. no. e2022JF006850.

[61] K. L. Riverman et al., "Is Lake Europa really a lake? Geophysical observations of a possible water body in the NW Greenland accumulation zone," in *Proc. AGU Fall Meeting Abstracts*, 2023, p. C33A.