High-Accuracy Classification of Radiation Waveforms of Lightning Return Strokes

Ting Wu¹, Daohong Wang¹, and Nobuyuki Takagi¹

¹Gifu University

February 20, 2023

Abstract

A machine-learning classifier for radiation waveforms of negative return strokes (RSs) is built and tested based on the Random Forest classifier using a large dataset consisting of 14,898 negative RSs and 159,277 intracloud (IC) pulses with 3-D location information. Eleven simple parameters including three parameters related with pulse characteristics and eight parameters related with the relative strength of pulses are defined to build the classifier. Two parameters for the evaluation of the classifier performance are also defined, including the classification accuracy, which is the percentage of true RSs in all classified RSs, and the identification efficiency, which is the percentage of correctly classified RSs in all true RSs. The tradeoff between the accuracy and the efficiency is examined and simple methods to tune the tradeoff are developed. The classifier achieved the best overall performance with an accuracy of 98.84% and an efficiency of 98.81%. With the same technique, the classifier for positive RSs is also built and tested using a dataset consisting of 8,700 positive RSs. The classifier has an accuracy of 99.04% and an efficiency of 98.37%. We also demonstrate that our classifiers can be readily used in various lightning location systems. By examining misclassified waveforms, we show evidence that some RSs and IC discharges produce special radiation waveforms that are almost impossible to correctly classify without 3-D location information, resulting in a fundamental difficulty to achieve very high accuracy and efficiency in the classification of lightning radiation waveforms.

High-Accuracy Classification of Radiation Waveforms of 1 Lightning Return Strokes 2

Ting Wu¹, Daohong Wang¹, Nobuyuki Takagi¹

¹Department of Electrical, Electronic and Computer Engineering, Gifu University, Gifu, Japan

Key Points: 5

3

4

6	• A machine-learning classifier for negative return strokes is built using a large dataset
7	with 3-D location information
8	- Both an accuracy and an efficiency of about 98.8% are achieved and the accuracy-efficiency
9	tradeoff can be easily controlled
10	• Some return strokes and IC discharges produce special waveforms that are fundamentally
11	difficult to classify without 3-D location results

difficult to classify without 3-D location results

Corresponding author: T. Wu, wu.ting.x4@f.gifu-u.ac.jp

12 Abstract

A machine-learning classifier for radiation waveforms of negative return strokes (RSs) 13 is built and tested based on the Random Forest classifier using a large dataset consisting 14 of 14,898 negative RSs and 159,277 intracloud (IC) pulses with 3-D location information. 15 Eleven simple parameters including three parameters related with pulse characteristics 16 and eight parameters related with the relative strength of pulses are defined to build the 17 classifier. Two parameters for the evaluation of the classifier performance are also defined, 18 including the classification accuracy, which is the percentage of true RSs in all classified 19 RSs, and the identification efficiency, which is the percentage of correctly classified RSs 20 in all true RSs. The tradeoff between the accuracy and the efficiency is examined and 21 simple methods to tune the tradeoff are developed. The classifier achieved the best overall 22 performance with an accuracy of 98.84% and an efficiency of 98.81%. With the same technique, 23 the classifier for positive RSs is also built and tested using a dataset consisting of 8,700 24 positive RSs. The classifier has an accuracy of 99.04% and an efficiency of 98.37%. We 25 also demonstrate that our classifiers can be readily used in various lightning location systems. 26 By examining misclassified waveforms, we show evidence that some RSs and IC discharges 27 produce special radiation waveforms that are almost impossible to correctly classify without 28 3-D location information, resulting in a fundamental difficulty to achieve very high accuracy 29

³⁰ and efficiency in the classification of lightning radiation waveforms.

³¹ Plain Language Summary

Lightning location systems are required to classify return strokes (RSs) from intracloud 32 discharges accurately and efficiently because the RS is the main discharge component 33 that poses direct threats to the human society. In this paper, we report a machine-learning 34 classifier for negative RSs built using a large dataset with accurate 3-D location information. 35 The classifier has an accuracy of 98.84% (98.84% of classified RSs are correct classifications) 36 and an efficiency of 98.81% (98.81% of RSs can be correctly classified). With the same 37 technique, we also built a classifier for positive RSs with similarly high accuracy and efficiency. Our classifiers only require some simple waveform parameters and can be readily used 39 in various national and continental lightning location systems. A sample Python script 40 to use the classifier is provided and readers are encouraged to test the classifier using their 41 own dataset. We also demonstrate that some RSs and intracloud discharges produce abnormal 42 waveforms, so 100% accuracy or efficiency is fundamentally difficult to realize using only 43 waveform information. 44

45 1 Introduction

Ground-based lightning location systems (LLSs) are widely used to monitor lightning 46 activities. A prominent feature of ground-based LLSs is that lightning activities in a wide 47 area can be monitored in real time with only a limited number of sensors. Some famous 48 national and continental LLSs include the National Lightning Detection Network (NLDN) 49 covering the continental United States (e.g. Cummins & Murphy, 2009), the European 50 Cooperation for Lightning Detection network (EUCLID) covering the European continent 51 (e.g. Schulz et al., 2016), and the Earth Networks Total Lightning Network (ENTLN) 52 (e.g. Zhu et al., 2022) with the aim of a global coverage. 53

It is a basic requirement for LLSs to automatically and efficiently classify cloud-to-ground (CG) lightning flashes from intracloud (IC) flashes as the former consist of discharges with direct connections to the ground and thus pose a much larger threat to the human society. The fundamental difference between a CG flash and an IC flash is that a CG flash contains one or more return strokes (RSs), so the classification of CG flashes is basically realized by classifying RSs. Further, it is well known that RSs produce characteristic electric field radiation waveforms that are largely different from those of IC discharges (e.g. Lin et al., 1979), so most LLSs classify RSs based on their waveform characteristics.

However, RSs actually can produce radiation waveforms with a variety of special 62 features under some special conditions. For example, some RSs in winter thunderstorms 63 are known to produce abnormal radiation waveforms, some of which could not be correctly 64 classified by LLSs (Wu, Wang, & Takagi, 2021; Wu, Wang, Huang, & Takagi, 2021). It is also well known that RSs striking tall objects produce much narrower radiation waveforms (Pavanello et al., 2007; Zhu et al., 2018). On the other hand, IC discharges include various 67 discharge processes such as narrow bipolar events and recoil leaders, some of which may 68 produce radiation waveforms with certain similar features as RS waveforms. As a result, for most LLSs, it is basically very difficult to achieve a very high classification accuracy 70 of RSs. For example, Zhu et al. (2016) reported that out of 339 RSs in Florida in 2014 71 that were also recorded by the NLDN, 312 (92%) were correctly classified as RSs by the 72 NLDN. Kohlmann et al. (2017) reported that the classification accuracy of EUCLID for 73 RSs were generally around 90% based on ground-truth data in various regions of Europe. 74 For some particular thunderstorms or some special types of discharges, misclassifications 75 by LLSs can be more common. For example, Fleenor et al. (2009) found that 204 out 76 of 376 (54%) of RSs reported by the NLDN during a field campaign in 2005 were actually 77 IC discharges. Leal et al. (2019) found that compact intracloud discharges with estimated 78 peak currents larger than 50 kA were all falsely classified as RSs by both NLDN and ENTLN. 79 Paul et al. (2020) reported that out of 40 RSs detected at the Peissenberg Tower, 12 (30%) were falsely classified as IC discharges. 81

In order to overcome the uncertainties in classifications based only on radiation waveforms, Betz et al. (2004) proposed a pseudo 3-D technique to assist the discrimination of RSs and IC discharges based on the fact that the elevation of IC discharges would have some contributions to the time delay. However, this technique also has some limitations. For example, IC discharges need to have significant elevations, the baseline of the LLS cannot be too long, and lightning discharges first need to be located accurately in 2-D. These limitations prevented the wide implementation of this technique.

In recent years, machine-learning techniques have been developing rapidly, and these techniques seem to be promising in significantly increasing the classification accuracy of lightning radiation waveforms. Wang et al. (2020) developed a convolutional neural network 91 to classify radiation waveforms of lightning discharges recorded by the Advanced Direction-time 92 Lightning Detection System in China. They reported an accuracy of over 99%. However, 93 they apparently did not have the height information of lightning discharges and thus could 94 not unambiguously differentiate RSs and IC discharges, so the accuracy remains questionable. 95 Zhu et al. (2021) used the Support Vector Machines (SVM) model to classify CG and 96 IC flashes recorded by the Cordoba Marx Meter Array. The lightning data were in 3-D, 97 so they could employ the discharge height information to build a dataset with accurate 98 discharge types. They reported an overall accuracy of 97%. However, their proposed method 99 requires full waveform information, while most LLSs only retrieve a few parameters of 100 electric field waveforms of lightning discharges, making it somewhat difficult for existing 101 systems to adopt the method. 102

In this paper, we report a simple yet high-accuracy machine-learning technique based 103 on the Random Forest classifier to classify RSs. We will use a large dataset containing 104 about 15,000 negative RSs and many more IC discharges with accurate 3-D location information 105 to train and test the classifier. As will be described in this paper, many of the recorded 106 RSs and IC discharges produced atypical radiation waveforms that were challenging to 107 be correctly classified. However, the accuracy of our classifier is close to 99% demonstrated 108 by evaluations in various respects. Our classifier requires only some simple parameters 109 of lightning radiation waveforms, so it can be readily used by most LLSs. 110



Figure 1. (a) Negative RSs (black dots) observed from July 19 to August 26 in 2017. (b) Positive RSs observed from September 26, 2021 to September 3, 2022. Red squares represent observation sites of FALMA.

¹¹¹ 2 Observation and Data

During the summer of 2017, we set up a low-frequency (LF) lightning mapping system 112 called Fast Antenna Lightning Mapping Array (FALMA) in central Japan. The FALMA 113 consisted of 12 sites covering an area of about 80×80 km². Locations of these 12 sites are shown as red squares in Figure 1a. At every site, a fast antenna working in the frequency 115 band of 500 Hz to 500 kHz was used to receive radiation signals from lightning discharges. 116 The signals were recorded with a sampling rate of 25 MS/s. As described by Wu et al. 117 (2018a), thanks to improvements made in both the hardware and the software, we realized 118 high-quality 3-D lightning mapping with the FALMA. As can be seen from examples of 119 lightning flashes in Wu et al. (2018a) and Wu et al. (2019), 3-D mapping results of FALMA 120 have similar quality to those of very-high-frequency (VHF) systems such as the Lightning 121 Mapping Array (Rison et al., 1999).

Data obtained from July 19 to August 26 are used in this study for building and 123 testing the classifier for negative RSs. All data are reprocessed for this study. The largest 124 positive pulse (the same polarity as the negative RS, using the atmospheric electricity 125 sign convention) in each 20-ms window is located in 3-D. Only discharges located in the region shown in Figure 1a, a $90 \times 90 \text{ km}^2$ area over the FALMA network, are used in order 127 to ensure reliable 3-D locating. Pulses with source heights lower than 500 m are treated 128 as candidates of RSs. Their waveforms are then confirmed manually, and for some ambiguous 129 pulses, they are further manually located to determine their source heights. In this way, 130 we can unambiguously determine that the selected pulses are truly RSs. The number 131 of IC discharges are much larger than that of RSs, so we cannot manually confirm waveforms 132 of all IC discharges, and we only use pulses with source heights larger than 3000 m as 133 IC pulses. There are 14,898 pulses confirmed as negative RSs and 159,277 pulses as IC discharges. Locations of these RSs are shown as black dots in Figure 1a. It should be 135 noted that we will build a classifier for negative RSs rather than negative CG flashes; 136 a CG flash consists of at least one RS and also many IC discharges, both of which need 137 to be correctly classified. 138

Using the high-quality dataset of 2017 summer, we will establish the technique for 139 building the classifier as will be described in Sections 3.1 to 3.5. Further, using the same 140 technique, we will also build a classifier for positive RSs as will be described in Section 3.6. 141 However, positive RSs in central Japan in summer are quite rare (Wu et al., 2018b). In order to accumulate a large number of positive RSs, we will use the data collected during 143 a long period, from September 26, 2021 to September 3, 2022. During this period, we 144 set up a FALMA network covering a large area for 2-D locating of both summer and winter 145 lightning. Observation sites are shown as red squares in Figure 1b. A total of 8700 positive 146 RSs observed in an area with a radius of 300 km are identified and will be used for building 147 and testing the classifier for positive RSs. Locations of these positive RSs are shown as 148 black dots in Figure 1b. The procedure for the identification of these positive RSs will 149 be further described in Section 3.6. 150

Our classifiers will be built and tested mainly based on the Random Forest classifier,
 which is one of the most widely used machine-learning models for classification tasks.
 A brief comparison will also be made with the SVM classifier, another popular machine-learning model, in Section 3.4.

- ¹⁵⁵ 3 Methods and Results
- 156

3.1 Method to Evaluate the Performace of a Classifier

Before building the classifier, first we need to define some parameters as indicators of the performance of a classifier. One obvious parameter to evaluate the performance is the classification accuracy, or simply *accuracy*, that is, the percentage of true RSs in the waveforms classified as RSs. However, only this parameter is apparently not enough, as it is always possible to build a classifier with very strict criteria so that it only identifies very typical RS waveforms. Another important parameter is the identification efficiency, or simply *efficiency*, that is, the percentage of correctly classified RSs in all RSs.

Suppose the number of RSs is N_R , and the number of IC discharges is N_I . Of the N_R RSs, N_{Rc} are correctly classified (the subscript *c* stands for "correct"), and the remaining $N_R - N_{Rc}$ are misclassified as IC discharges. Of the N_I IC discharges, N_{Ic} are correctly classified, and the remaining $N_I - N_{Ic}$ are misclassified as RSs. The accuracy and the efficiency are defined as follows.

$$Accuracy = \frac{N_{Rc}}{N_{Rc} + (N_I - N_{Ic})} \tag{1}$$

169

$$Efficiency = \frac{N_{Rc}}{N_R} \tag{2}$$

During the process to build the classifier, we will experiment and tune various parameters of the classifier to make the accuracy and the efficiency as high as possible.

Normally a dataset is split into a larger training set and a smaller test set, with 172 the training set used to train a classifier and the test set used to test or evaluate the performance 173 of the classifier. In this study, we use an improved approach. All RS and IC data are combined, 174 shuffled and then divided into five equal parts. Each part is in turn used as the test set 175 and the remaining four parts combined are used as the training set. In this way, a classifier 176 is built and tested for five times and five results of accuracy and efficiency are calculated. The average values of five tests will be used as the final results. In this way, we can avoid 178 any random biases in the test set. Moreover, as will be described in Section 4, in this 179 way all data can be tested and we can find as many atypical waveforms as possible that 180 are difficult to be correctly classified. 181

¹⁸² 3.2 Waveform Paramaterization

We will define some waveform parameters to be used for building the classifier. First we describe the procedure to calculate waveform parameters based on multiple-site records. As waveforms recorded at a close distance contain the electrostatic and induction field components (e.g. Thottappillil et al., 1997) that may significantly distort the waveforms, observation sites within 40 km from a discharge are first excluded. Waveforms recorded by the remaining sites are used to calculate the parameters, and for each parameter, the median value of the results calculated based on these sites are used as the final result of the parameter for the discharge.

191

3.2.1 Parameters Related with Pulse Characteristics

First we define three basic parameters related with pulse characteristics. Definitions of these parameters are illustrated using an RS pulse in Figure 2a and an IC pulse in Figure 2b (blue parameters).

- 195 1. T_{rise} : The rise time of a pulse (10% to peak).
- **196** 2. T_{fall} : The fall time of a pulse (peak to zero).
- **3.** T_{half} : The pulse width at the half maximum.

With only these three basic parameters, we trained and tested the Random Forest 198 classifier using the negative RS and IC dataset obtained in 2017 summer. As described 199 in Section 3.1, the dataset is divided into five parts and each part in turn is used as the 200 test set, so the classifier is trained and tested for five times. The accuracy ranges from 201 72.25% to 73.57% with an average of 72.82%, and the efficiency ranges from 70.80% to 72.81% with an average of 71.59%. We also tried to add two related parameters, including 203 the pulse width, which is the sum of the rise time and fall time, and the ratio of fall time 204 to rise time, but the result has little difference (the average accuracy is 72.17% and the 205 average efficiency is 70.86%). 206

207 Indeed, with only these basic pulse parameters, it is difficult to accurately classify208 RSs.

209

3.2.2 Parameters Related with Relative Strength

An important feature of the RS waveform is that pulses right before and after an RS pulse is usually much weaker. The following parameters are defined to employ this feature. These parameters are also illustrated in Figures 2a and 2b.

- 1. R_{bp1} : The ratio of A_0 to A_{bp1} , in which A_0 is the peak amplitude of the target pulse, and A_{bp1} is the maximum amplitude of pulses right before the target pulse (from -100 μ s to 10% peak) as illustrated in Figure 2. The subscript *b* stands for "before", and the subscript *p* stands for "positive".
- 217 2. R_{bn1} , R_{bp2} , R_{an1} , R_{an1} , R_{ap2} , R_{an2} : These parameters are defined in the 218 same way as R_{bp1} , also illustrated in Figure 2. Note that the subscript *a* stands 219 for "after", and the subscript *n* stands for "negative".

The three parameters defined in Section 3.2.1 along with the eight new parameters defined above are used to train the Random Forest classifier. The accuracy of five tests ranges from 98.86% to 99.32% with an average of 99.02%, and the efficiency ranges from 98.02% to 98.66% with an average of 98.34%. It is clear that these new parameters representing the relative strength are very effective in the classification of RSs.



Figure 2. Illustration of waveform parameters using (a) an RS pulse and (b) an IC pulse. (c) Relative importance of waveform parameters.

3.2.3 Parameter Importance

The Random Forest classifier outputs a value indicating the relative importance 226 of each parameter in contributing to the performance, from which we can evaluate the 227 effectiveness of each parameter in the classification of RSs. The results are shown in Figure 2c. 228 Values of the importance of all parameters combined equal to 1. We can see that parameters 229 related with the pulse strength relative to previous pulses (red parameters in Figure 2) 230 are generally more important than other parameters. This is easy to understand as an 231 RS pulse is preceded by leader pulses which are usually much weaker than the RS pulse. 232 By contrary, an IC pulse is usually preceded by other IC pulses with comparable amplitudes. Therefore, parameters related with the relative strength are very effective in the classification 234 of RSs. 235

We can also see that parameters related with pulse characteristics (blue parameters) have relatively low importance, which is why the classifier performance is very poor with only these parameters as described in Section 3.2.1. It also indicates that traditional RS classification methods based on pulse characteristics are not very reliable.

3.3 Tradeoff Between Accuracy and Efficiency

From the above result, we can see one feature of the classifier is that the accuracy is always higher than the efficiency. It is obvious that increasing the efficiency usually implies decreasing the accuracy. However, it is desirable if we can control the tradeoff between the accuracy and the efficiency. For example, in some situations, it may be required to identify as many RSs as possible, so a high efficiency is essential while a low accuracy is tolerable. Next we will investigate two factors that influence the tradeoff between the accuracy and the efficiency.

248

240

3.3.1 Influence of Sample Size Imbalance

One reason for the higher accuracy in the classifier built in the previous section is a much larger sample of IC discharges compared with the sample of RSs. With such a biased dataset, the classifier is more likely to misclassify RSs, as also noted by Zhu et al. (2021). We can simply duplicate the sample of RSs to make the classifier identify more RSs, though at the cost of more misclassifications of IC discharges. Note that the duplication should only be made for the training set.

With the original dataset, 247 of 14,898 RSs (1.7%) are misclassified, but only 145 255 of 159,277 IC pulses (0.091%) are misclassified. If we duplicate the dataset of RSs in the 256 training set, the number of misclassified IC pulses increases to 162 while the number of 257 misclassified RSs decreases to 199. We tried to make more duplications and tested the 258 classifier, and the results of the accuracy and the efficiency are shown in Figure 3a. With 259 one duplication of the RS training set, the accuracy decreases from 99.02% to 98.91%but the efficiency increases from 98.34% to 98.66%. With two duplications, the accuracy decreases to 98.84% but the efficiency increases to 98.76%, very close to the accuracy. 262 With further duplications, we can see that both the accuracy and the efficiency are generally 263 very similar, changing between 98.75% and 98.85%, indicating that the sample size imbalance 264 does not have a significant effect any more. 265

If we use the average value of the accuracy and the efficiency as the indicator of the overall performance of a classifier, we can see from Figure 3a that with four duplications of the RS training set, the classifier has the highest performance with an accuracy of 98.84% and an efficiency of 98.81%. We treat this as the best performance of the classifier for negative RSs and this classifier will be used for further evaluations in the following section.



Figure 3. Variations of the accuracy and the efficiency with (a) different times of duplications of RS training set and (b) different thresholds of probability to classify RSs.

	Eleven Parameters (Section 3.2.2)		
Classifier	Accuracy (%)	Efficiency (%)	Time Cost (seconds)
Random Forest	99.02	98.34	20
SVM	98.43	97.42	96
Du	plicating RS Tra	ining Set (Section	n 3.3.1)
Classifier	Accuracy (%)	Efficiency (%)	Time Cost (seconds)
Random Forest	98.84	98.81	28
SVM	97.98	98.09	108

 Table 1. Comparison of the Random Forest classifier and the SVM classifier

271 3.3.2 Influence of Probability Thresholds

When classifying a pulse, the Random Forest classifier can output the probability
that the pulse is a true RS. By default, the classifier determines a pulse as an RS when
the probability is larger than 50%. By changing the probability threshold, we can conveniently
tune the accuracy-efficiency tradeoff.

Figure 3b shows variations of the accuracy and the efficiency related with the probability threshold. We can see that as the probability threshold increases, the accuracy increases while the efficiency decreases. This is easy to understand; a higher probability threshold represents stricter criteria to classify RSs, so naturally the identified RSs are more likely true RSs (higher accuracy), but at the same time fewer RSs can be identified (lower efficiency). In practice, when using the classifier we can set a customized probability threshold that fits the specific requirements of an application to achieve desired accuracy or efficiency.

283

3.4 Comparison of Different Machine-learning Models

Apart from the Random Forest classifier, another popular machine-learning model 284 for classification is the SVM classifier, which was used by Zhu et al. (2021) for the classification 70E of lightning pulses. Here we make a brief comparison of the Random Forest and the SVM classifiers. First we use the scheme described in Section 3.2.2 (using 11 parameters illustrated 287 in Figure 2) to train the classifiers, and the results are shown in Table 3.4 (upper part). 288 We can see the SVM classifier has slightly lower accuracy and efficiency than the Random 289 Forest classifier. Further, we use the scheme described in Section 3.3.1 (duplicating the 290 RS training dataset) to train the classifiers, and again, the SVM classifier has slightly 291 lower accuracy and efficiency. Another difference is in the time needed to train a classifier; 292 it takes less than 30 seconds to train an Random Forest classifier while the time needed 293 to train an SVM classifier is around 100 seconds. A significantly shorter time to build a classifier is potentially very useful as it would be more convenient to experiment various 295 combinations of parameters in order to boost the performance of the classifier. 296

297

3.5 Testing Using Remote Lightning Discharges

Lightning discharges used for training and evaluating classifiers described above are all very close to most of FALMA sites in order to ensure the 3-D location accuracy. However, many LLSs, especially national and continental LLSs, have long baselines of a few hundred kilometers, so lightning discharges observed by these systems are generally very far away from most observation sites. Therefore, it is desirable to evaluate the performance of a classifier for remote lightning discharges. We use lightning discharges located more than 150 km away from the center of the FALMA network in 2017 summer (the origin in Figure 1a) for this investigation. At such a large distance, only a small number of discharges can be located with sufficient accuracy, and we can only make 2-D locating, so we cannot classify RSs using the height information. Therefore, we manually inspected waveforms of all located events and determine their types.

There are a total of 594 located pulses. The classifier described in Section 3.3.1 (the 310 training set duplicated for four times) are used to classify these pulses. A total of 361 311 pulses were classified as RSs, and there was no clear misclassification. The remaining 233 312 pulses were classified as IC discharges, and four of them were likely RSs. However, it should 313 be noted that as there is no height information for these pulses, it is sometimes difficult 314 to determine the true discharge type, so it is possible that there were actually more misclassifications. 315 Assuming there are only four RSs misclassified as IC discharges, from Equations 1 and 316 2, we can get an accuracy of 100% and an efficiency of 98.9%. Note that when detecting 317 remote lightning discharges, as in the case of long-baseline LLSs, only a small portion 318 of IC discharges that are relatively strong can be located, so the chance of misclassifying an IC pulse as an RS is relatively low, which may be one reason for the 100% accuracy 320 in this evaluation. 321

The above results demonstrated that our classifier also has good performance when classifying remote RSs, so the classifier can also be used in long-baseline LLSs.

324 3.6 Classification of Positive Return Strokes

The methods described above can also be used to build a classifier for the classification 325 of positive RSs. However, positive CG flashes are very rare in summer thunderstorms 326 in central Japan. As reported by Wu et al. (2018b), only 46 positive CG flashes consisting 327 of 53 positive RSs were observed and could be located in 3-D during the summer observation 328 of 2017. Therefore, here we also include the data obtained in other periods. First we use the 690 positive RSs observed during the winter of 2018 (Wu et al., 2022) to build a preliminary 330 classifier for the identification of positive RSs. Then we use this classifier to search the 331 data recorded in about one year from September of 2021 for possible positive RSs. As 332 described in Section 2, during this period, we set up a FALMA network with long baselines 333 for 2-D locating of both summer and winter lightning. Waveforms of the identified positive 334 RSs by the preliminary classifier are manually confirmed to exclude obvious false classifications. 335 Indeed, the preliminary classifier identified many pulses that were clearly IC pulses and 336 we painstakingly excluded all apparent IC pulses by manual inspections. In this way, we 337 collected the data of 8700 positive RSs, locations of which are shown in Figure 1b. Note 338 that there is no height information for these positive RSs, so this dataset is not as accurate 339 as the negative RS dataset in 2017 summer used in previous sections. 340

For IC data, we also use the data of summer observation of 2017 as these data have 341 accurate 3-D location results. However, different from the IC dataset for the negative 342 RS classifier, IC pulses for building positive RS classifier should have the same polarity 343 as positive RSs. So we located IC pulses having the same polarity as positive RSs and 344 selected those with heights larger than 3 km, the same treatment as that in building the 345 negative RS classifier. On the other hand, as the size of positive RS dataset is relatively 346 small, we do not need too many IC data, so for simplicity, we only located one IC pulse 347 in every 50-ms window. Finally, we collected a total of 113,922 IC pulses. 348

Using these datasets, and with the same scheme for building negative RS classifier described in Section 3.3.1, we built and tested the classifier for positive RSs. It is found that with the RS training set duplicated for one time, the classifier has the best overall performance. It has an accuracy of 99.04% and an efficiency of 98.37%, generally similar to the performance of the negative RS classifier. This result demonstrated that as long as there are enough data of positive RSs, we can also build a high-accuracy classifier for positive RSs in the same way as building the negative RS classifier.

Although the dataset of positive RSs does not have 3-D location information and thus is not as accurate as the negative RS dataset, as positive RSs are much rarer than negative RSs and it is very difficult to collect a large and reliable sample, we believe our classifier is very valuable for future observations and researches. Moreover, as all waveforms of identified positive RSs have been manually confirmed, the classifier likely has an accuracy similar to that of the manual classification.

³⁶² 4 Atypical Intracloud and Return Stroke Waveforms

As described in Section 3.1, the whole dataset of 2017 summer is divided into five parts, with each part in turn used as the testing set and the remaining four parts combined used as the training set. In this way, all pulses can be tested and we can identify as many pulses as possible that are potentially difficult to classify. Using the classifier built in Section 3.3.1 with the RS training set duplicated for four times, all pulses are classified. Of the 14,898 RSs, 178 were misclassified as IC discharges, and of the 159,277 IC pulses, 173 were misclassified as RSs. Waveform figures of all these misclassified pulses are provided in the data repository.

There are several common reasons for misclassifications of RS pulses as IC pulses. 370 Waveforms of four examples are shown in Figures 4a-d, all of which are misclassified as 371 IC discharges and whose source heights have been confirmed to be close to the ground. 372 First, it is well known that RSs striking tall grounded objects usually produce very narrow 373 pulses (Araki et al., 2018; Cai et al., 2022; Pavanello et al., 2007; Zhu et al., 2018), making 374 it easy to misclassify them as IC discharges. One example is shown in Figure 4a. This pulse is located near a transmission tower, and its pulse width is only about 6 μ s, indicating 376 that it is likely produced by an RS striking the tower. Second, two RSs sometimes occur 377 sequentially with a very small time difference of a few tens of microseconds, and if the 378 second RS has a larger peak than the first one, the second RS may be misclassified as 379 an IC pulse. One example is shown in Figure 4b. Such RSs are likely the so-called "multiple-termination 380 strokes" (Kong et al., 2009; Sun et al., 2016) or "forked strokes" (Ballarotti et al., 2005), 381 with two RSs induced by two branches of the same leader. Third, an RS may occur almost simultaneously with IC discharges of other lightning flashes, resulting in a peculiar waveform 383 and thus misclassified as an IC discharge. One example is shown in Figure 4c. While the 384 positive pulse is confirmed to be produced by an RS, the two negative pulses labeled as 385 "IC" are produced by IC discharges in an independent lightning flash and are located about 386 87 km away from the RS. The resultant waveform appears to be abnormal and is difficult 387 to be identified as an RS. Finally, some RSs apparently produce waveforms that are largely 388 different from typical RS waveforms but the reason is not yet clear. One example is shown in Figure 4d. The pulse has a rise time of about 18 μ s while its fall time is only about 9 μs . 391

Another example of an RS producing abnormal waveform is shown in Figure 5 along with location results of the preceding leader. This RS is a subsequent RS. We can see from the location results in Figure 5a a dart leader with a speed of about 4×10^6 m/s preceding the RS, and the RS is located very close to the ground as indicated by the cross sign. From Figure 5c, we can see details of the RS waveform. It contains two peaks with the second peak much larger than the first one, resulting in a much larger rise time than the fall time. Without the 3-D location results, it is very difficult to determine that the waveform is produced by an RS.

The major reason for IC discharges misclassified as RSs is that waveforms of some
IC discharges have some similar features as those of RSs. Four examples are shown in
Figures 4e-h. All of these waveforms appear very similar to those of RSs. However, their
source heights range from 5.9 to 15.1 km, indicating that they are produced by IC discharges.



Figure 4. (a)-(d) Atypical E-change waveforms produced by RSs but misclassified as IC discharges. (e)-(h) Atypical E-change waveforms produced by IC discharges but misclassified as RSs. The value of d represents the distance between the discharge and the observation site recording the waveform. The value of h represents the source height of the IC discharge.



Figure 5. Location result and E-change waveforms of an RS misclassified as an IC discharge. (a) Height-time location results of the dart leader preceding the RS. The cross sign represents the RS. (b) E-change waveform of the RS and preceding discharges. (c) E-change waveform of the RS. The value of d represents the distance between the RS and the observation site recording the waveform.

We also manually located these pulses to make sure that there were no large errors in the source height results. We can see that these pulses have relatively short rise times and much longer fall times. Pulses in Figures 4e-g also have fine structures superimposed on the falling part, similar to waveforms of first RSs, and the pulse in Figure 4h resembles the waveform of a subsequent RS. These similar features as RS waveforms make it almost impossible to correctly classify them as IC discharges without the 3-D location results.

Another example of an IC pulse appearing similar to RS pulses is shown in Figure 6 410 along with location results of preceding discharges. From the height-time location results in Figure 6a, we can see that a leader first propagated above 6.5 km and then descended 412 to a height of about 5.5 km, and then the large IC pulse is produced, represented by the 413 cross sign. From the E-change waveform in Figure 6c, we can see that the large IC pulse 414 is very similar to an RS pulse, with preceding pulses resembling stepped leader pulses. 415 With the help of the 3-D location results, we can be sure that this RS-like pulse is produced 416 by IC discharges. We are not aware of any study reporting such RS-like IC pulses. In 417 our future studies, we will explore the mechanism responsible for these special IC pulses. 418

These examples of special RS and IC waveforms illustrate the fact that some RSs and IC discharges produce atypical radiation waveforms from which the discharge types cannot be accurately determined, resulting in a fundamental difficulty to achieve very high accuracy and efficiency using only waveform information. This result also illustrates the importance of accurate 3-D location results in scientific investigations of lightning phenomena.

425 5 Conclusions

Using a large dataset with 3-D location results, we built a classifier for radiation 426 waveforms of negative RSs based on the Random Forest classifier. Eleven simple parameters 427 are defined for building the classifier, including three parameters related with pulse characteristics 128 and eight parameters related with relative strength of pulses. A classification accuracy 429 of 98.84% and an identification efficiency of 98.81% are achieved. We also demonstrated 430 methods to tune the tradeoff between the accuracy and the efficiency so the classifier can 431 be used in applications with different requirements of the accuracy or the efficiency. Although 432 the classifier is built based on the observation of a compact lightning mapping system, 433 we demonstrated that the classifier also has high accuracy and efficiency for remote lightning 434 discharges and can be readily used in long-baseline LLSs. With the same methods, we 435 also built a classifier for positive RSs which has similarly high accuracy and efficiency as the classifier for negative RSs. 437

Misclassified RS and IC waveforms are examined and some common reasons for misclassifications
are analyzed. We demonstrated that RSs sometimes produce radiation waveforms that
are largely different from normal RS waveforms, and IC discharges sometimes produce
waveforms that appear very similar to RS waveforms. Therefore, some RS and IC waveforms
are fundamentally difficult to be correctly classified without 3-D location information,
and it is likely that such misclassifications commonly exist in most LLSs. The results
also imply the importance of 3-D location results in detailed analyses of lightning phenomena.

445 Open Research Section

Datasets for building and testing the classifiers as well as waveform figures of all positive and negative RSs can be found at https://doi.org/10.5281/zenodo.7641792. Sample Python scripts for using the classifiers will be made publicly available after the acceptance of this paper.



Figure 6. Location result and E-change waveforms of an IC pulse misclassified as an RS pulse. (a) Height-time location results of the IC pulse and preceding discharges. The cross sign represents the location of the IC pulse. (b) E-change waveform of the IC pulse and preceding discharges. (c) E-change waveform of the IC pulse. The value of *d* represents the distance between the IC discharge and the observation site recording the waveform.

450 Acknowledgments

This study was supported by the Ministry of Education, Culture, Sports, Science, and Technology of Japan (Grants 20H02129 and 21K03681).

453 References

454	Araki, S., Nasu, Y., Baba, Y., Rakov, V. A., Saito, M., & Miki, T. (2018, sep).
455	3-d finite difference time domain simulation of lightning strikes to the 634-m
456	tokyo skytree. Geophysical Research Letters, $45(17)$, $9267-9274$. doi:
457	$\frac{1}{10000000000000000000000000000000000$
458	observations of negative ground flashes on a millisecond-scale <i>Ceonhysical</i>
459	Research Lettere 32(23) doi: https://doi.org/10.1020/2005gl023880
460	Betz H D Schmidt K Oettinger P $\&$ Wirz M (2004 jun) Lightning
401	detection with 3-d discrimination of intracloud and cloud-to-ground discharges
402	Geophysical Research Letters 31(11) n/a-n/a doi: 10.1029/2004gl019821
464	Cai, L., Liu, W., Zhou, M., Wang, J., Yan, R., Tian, R., & Fan, Y. (2022, dec).
465	Differences of electric field parameters for lightning strikes on tall towers and
466	nonelevated objects. <i>IEEE Transactions on Electromagnetic Compatibility</i> ,
467	64(6), 2113–2121. doi: https://doi.org/10.1109/temc.2022.3207237
468	Cummins, K. L., & Murphy, M. J. (2009, aug). An overview of lightning locating
469	systems: History, techniques, and data uses, with an in-depth look at the u.s.
470	NLDN. IEEE Transactions on Electromagnetic Compatibility, 51(3), 499–518.
471	doi: https://doi.org/10.1109/temc.2009.2023450
472	Fleenor, S. A., Biagi, C. J., Cummins, K. L., Krider, E. P., & Shao, XM.
473	(2009, feb). Characteristics of cloud-to-ground lightning in warm-season
474	thunderstorms in the central great plains. Atmospheric Research, 91(2-4),
475	333–352. doi: https://doi.org/10.1016/j.atmosres.2008.08.011
476	Kohlmann, H., Schulz, W., & Pedeboy, S. (2017, oct). Evaluation of EUCLID
477	IC/CG classification performance based on ground-truth data. In 2017 intermetional summasium on lightning metaction (VIV CIDDA) IEEE doi:
478	https://doi.org/10.1100/sindo.2017.8116806
479	Kong X Oio X Zhao V & Zhang T (2000 feb) Characteristics of
480	negative lightning flashes presenting multiple-ground terminations on a
482	millisecond-scale. Atmospheric Research, 91(2-4), 381–386. doi: https://
483	doi.org/10.1016/j.atmosres.2008.03.025
484	Leal, A. F., Rakov, V. A., & Rocha, B. R. (2019, aug). Compact intracloud
485	discharges: New classification of field waveforms and identification by lightning
486	locating systems. Electric Power Systems Research, 173, 251–262. doi:
487	https://doi.org/10.1016/j.epsr.2019.04.016
488	Lin, Y. T., Uman, M. A., Tiller, J. A., Brantley, R. D., Beasley, W. H., Krider,
489	E. P., & Weidman, C. D. (1979). Characterization of lightning return stroke
490	electric and magnetic fields from simultaneous two-station measurements.
491	Journal of Geophysical Research, 84 (C10), 6307. doi: https://doi.org/10.1029/
492	jc084ic10p06307
493	Paul, C., Heidler, F. H., & Schulz, W. (2020, feb). Performance of the european
494	lightning detection network EUCLID in case of various types of current
495	purses from upward fighting measured at the persenderg tower. IEEE Transactions on Floatromagnetic Commetibility $60(1)$ 116 122
496	$\frac{1100-123}{100} \text{ doi } 00000000000000000000000000000000000$
497	Pavanello D. Rachidi F. Janischewsky W. Rubinstein M. Hussein A. M.
498	Petrache E Jaquier A (2007 jul) On return stroke currents and remote
₩ ₽ ₽	electromagnetic fields associated with lightning strikes to tall structures.
501	2. experiment and model validation. Journal of Geophysical Research:
502	Atmospheres, 112(D13). doi: https://doi.org/10.1029/2006id007959

503	Rison, W., Thomas, R. J., Krehbiel, P. R., Hamlin, T., & Harlin, J. (1999, dec). A
504	GPS-based three-dimensional lightning mapping system: Initial observations
505	in central new mexico. <i>Geophysical Research Letters</i> , 26(23), 3573–3576. doi:
506	https://doi.org/10.1029/1999gl010856
507	Schulz, W., Diendorfer, G., Pedeboy, S., & Poelman, D. R. (2016, mar). The
508	european lightning location system EUCLID – part 1: Performance analysis
509	and validation. Natural Hazards and Earth System Sciences, 16(2), 595–605.
510	doi: https://doi.org/10.5194/nhess-16-595-2016
511	Sun, Z., Qie, X., Liu, M., Jiang, R., Wang, Z., & Zhang, H. (2016, jan).
512	Characteristics of a negative lightning with multiple-ground terminations
513	observed by a VHF lightning location system. Journal of Geophysical
514	Research: Atmospheres, 121(1), 413–426. doi: https://doi.org/10.1002/
515	2015id023702
516	Thottappillil, R., Bakov, V. A., & Uman, M. A. (1997, mar). Distribution of charge
517	along the lightning channel: Relation to remote electric and magnetic fields
E10	and to return-stroke models Journal of Geophysical Research: Atmospheres
510	102(D6) 6987–7006 doi: https://doi.org/10.1029/96id03344
519	Wang I Huang O Ma O Chang S He I Wang H Gao C (2020 feb)
520	Classification of VLF/LF lightning signals using sensors and deen learning
521	methods. Sensore $20(A)$ 1030 doi: https://doi.org/10.3300/s200/1030
522	Wu T Wang D Huang H & Takagi N (2021 nov) The strongest negative
523	lightning strokes in winter thunderstorms in ispan <i>Combusical Research</i>
524	Lettere /8(21) doi: https://doi.org/10.1020/2021gl005525
525	$M_{\rm H}$ T Wang D & Takagi N (2018a apr) Lightning mapping with an array of
520	fast antonnas Coonducical Research Letters /5(8) 3608-3705 doi: https://
527	doi org/10.1002/2018cl077628
528	$W_{\rm H}$ T Wang D & Takagi N (2018b aug) Logating proliminary breakdown
529	nulses in positive cloud to ground lightning — <i>Journal of Coophysical Research</i> :
530	Atmospheres, doi: https://doi.org/10.1020/2018id028716
531	$W_{\rm H}$ T Wang D fr Takagi N (2010 con) Valagities of positive leaders
532	in intracloud and nogative cloud to ground lightning flaches
533	of Coonhusical Research: Atmospheres 19/(17.18) 0083-0005
534	https://doi.org/10.1020/2010id030783
535	Wu T Wang D & Takagi N (2021 aug) Compact lightning strokes in winter
530	thunderstorms Journal of Geonhusical Research: Atmospheres 196(15) doi:
537	https://doi.org/10.1020/2021id034032
538	Wu T Wang D ℓ_r Takagi N (2022 nov) On the intensity of first return
539	strokes in positive cloud to ground lightning in winter — <i>Learnal of Coonhusical</i>
540	Beasearch: Atmospheres, 107(22), doi: https://doi.org/10.1020/2022id037282
541	The V Bitzer P Bakey V k Ding 7 (2021 jap) A machine learning approach
542	to classify cloud to ground and intracloud lightning <i>Coophwical Research</i>
543	Lettere $/8(1)$ doi: https://doi.org/10.1020/2020gl001148
544	The V Bakov V A Tran M D Lyu W & Micu D D (2018 con) Λ
545	modeling study of narrow electric field signatures produced by lightning strikes
546	to tall towars Journal of Coonhusical Research: Atmospheree, 199(18) doi:
547	https://doi.org/10.1020/2018id028016
548	The V Balow V A Tran M D & Nag A (2016 dec) A study of national
549	lightning detection network responses to network lightning based on ground
550	truth data acquired at LOC with emphasis on cloud discharge activity
551	Lowrnal of Coophysical Research, Atmospheres, 101(94), 14,651, 14,660
552	bttps://doi.org/10.1002/2016id025574
553	Thu V Stock M Lapione I & DiCanci F (2022 may) Upgrades of the conth
554	notworks total lightning network in 2021 — <i>Demote Considered 11(0)</i> 2200 — Join
555	networks total lightling network in 2021. <i>Remote Sensing</i> , $14(9)$, 2209. doi: https://doi.org/10.2200/rs14002200
556	nttps://doi.org/10.3390/1814092209

High-Accuracy Classification of Radiation Waveforms of 1 Lightning Return Strokes 2

Ting Wu¹, Daohong Wang¹, Nobuyuki Takagi¹

¹Department of Electrical, Electronic and Computer Engineering, Gifu University, Gifu, Japan

Key Points: 5

3

4

6	• A machine-learning classifier for negative return strokes is built using a large dataset
7	with 3-D location information
8	- Both an accuracy and an efficiency of about 98.8% are achieved and the accuracy-efficiency
9	tradeoff can be easily controlled
10	• Some return strokes and IC discharges produce special waveforms that are fundamentally
11	difficult to classify without 3-D location results

difficult to classify without 3-D location results

Corresponding author: T. Wu, wu.ting.x4@f.gifu-u.ac.jp

12 Abstract

A machine-learning classifier for radiation waveforms of negative return strokes (RSs) 13 is built and tested based on the Random Forest classifier using a large dataset consisting 14 of 14,898 negative RSs and 159,277 intracloud (IC) pulses with 3-D location information. 15 Eleven simple parameters including three parameters related with pulse characteristics 16 and eight parameters related with the relative strength of pulses are defined to build the 17 classifier. Two parameters for the evaluation of the classifier performance are also defined, 18 including the classification accuracy, which is the percentage of true RSs in all classified 19 RSs, and the identification efficiency, which is the percentage of correctly classified RSs 20 in all true RSs. The tradeoff between the accuracy and the efficiency is examined and 21 simple methods to tune the tradeoff are developed. The classifier achieved the best overall 22 performance with an accuracy of 98.84% and an efficiency of 98.81%. With the same technique, 23 the classifier for positive RSs is also built and tested using a dataset consisting of 8,700 24 positive RSs. The classifier has an accuracy of 99.04% and an efficiency of 98.37%. We 25 also demonstrate that our classifiers can be readily used in various lightning location systems. 26 By examining misclassified waveforms, we show evidence that some RSs and IC discharges 27 produce special radiation waveforms that are almost impossible to correctly classify without 28 3-D location information, resulting in a fundamental difficulty to achieve very high accuracy 29

³⁰ and efficiency in the classification of lightning radiation waveforms.

³¹ Plain Language Summary

Lightning location systems are required to classify return strokes (RSs) from intracloud 32 discharges accurately and efficiently because the RS is the main discharge component 33 that poses direct threats to the human society. In this paper, we report a machine-learning 34 classifier for negative RSs built using a large dataset with accurate 3-D location information. 35 The classifier has an accuracy of 98.84% (98.84% of classified RSs are correct classifications) 36 and an efficiency of 98.81% (98.81% of RSs can be correctly classified). With the same 37 technique, we also built a classifier for positive RSs with similarly high accuracy and efficiency. Our classifiers only require some simple waveform parameters and can be readily used 39 in various national and continental lightning location systems. A sample Python script 40 to use the classifier is provided and readers are encouraged to test the classifier using their 41 own dataset. We also demonstrate that some RSs and intracloud discharges produce abnormal 42 waveforms, so 100% accuracy or efficiency is fundamentally difficult to realize using only 43 waveform information. 44

45 1 Introduction

Ground-based lightning location systems (LLSs) are widely used to monitor lightning 46 activities. A prominent feature of ground-based LLSs is that lightning activities in a wide 47 area can be monitored in real time with only a limited number of sensors. Some famous 48 national and continental LLSs include the National Lightning Detection Network (NLDN) 49 covering the continental United States (e.g. Cummins & Murphy, 2009), the European 50 Cooperation for Lightning Detection network (EUCLID) covering the European continent 51 (e.g. Schulz et al., 2016), and the Earth Networks Total Lightning Network (ENTLN) 52 (e.g. Zhu et al., 2022) with the aim of a global coverage. 53

It is a basic requirement for LLSs to automatically and efficiently classify cloud-to-ground (CG) lightning flashes from intracloud (IC) flashes as the former consist of discharges with direct connections to the ground and thus pose a much larger threat to the human society. The fundamental difference between a CG flash and an IC flash is that a CG flash contains one or more return strokes (RSs), so the classification of CG flashes is basically realized by classifying RSs. Further, it is well known that RSs produce characteristic electric field radiation waveforms that are largely different from those of IC discharges (e.g. Lin et al., 1979), so most LLSs classify RSs based on their waveform characteristics.

However, RSs actually can produce radiation waveforms with a variety of special 62 features under some special conditions. For example, some RSs in winter thunderstorms 63 are known to produce abnormal radiation waveforms, some of which could not be correctly 64 classified by LLSs (Wu, Wang, & Takagi, 2021; Wu, Wang, Huang, & Takagi, 2021). It is also well known that RSs striking tall objects produce much narrower radiation waveforms (Pavanello et al., 2007; Zhu et al., 2018). On the other hand, IC discharges include various 67 discharge processes such as narrow bipolar events and recoil leaders, some of which may 68 produce radiation waveforms with certain similar features as RS waveforms. As a result, for most LLSs, it is basically very difficult to achieve a very high classification accuracy 70 of RSs. For example, Zhu et al. (2016) reported that out of 339 RSs in Florida in 2014 71 that were also recorded by the NLDN, 312 (92%) were correctly classified as RSs by the 72 NLDN. Kohlmann et al. (2017) reported that the classification accuracy of EUCLID for 73 RSs were generally around 90% based on ground-truth data in various regions of Europe. 74 For some particular thunderstorms or some special types of discharges, misclassifications 75 by LLSs can be more common. For example, Fleenor et al. (2009) found that 204 out 76 of 376 (54%) of RSs reported by the NLDN during a field campaign in 2005 were actually 77 IC discharges. Leal et al. (2019) found that compact intracloud discharges with estimated 78 peak currents larger than 50 kA were all falsely classified as RSs by both NLDN and ENTLN. 79 Paul et al. (2020) reported that out of 40 RSs detected at the Peissenberg Tower, 12 (30%) were falsely classified as IC discharges. 81

In order to overcome the uncertainties in classifications based only on radiation waveforms, Betz et al. (2004) proposed a pseudo 3-D technique to assist the discrimination of RSs and IC discharges based on the fact that the elevation of IC discharges would have some contributions to the time delay. However, this technique also has some limitations. For example, IC discharges need to have significant elevations, the baseline of the LLS cannot be too long, and lightning discharges first need to be located accurately in 2-D. These limitations prevented the wide implementation of this technique.

In recent years, machine-learning techniques have been developing rapidly, and these techniques seem to be promising in significantly increasing the classification accuracy of lightning radiation waveforms. Wang et al. (2020) developed a convolutional neural network 91 to classify radiation waveforms of lightning discharges recorded by the Advanced Direction-time 92 Lightning Detection System in China. They reported an accuracy of over 99%. However, 93 they apparently did not have the height information of lightning discharges and thus could 94 not unambiguously differentiate RSs and IC discharges, so the accuracy remains questionable. 95 Zhu et al. (2021) used the Support Vector Machines (SVM) model to classify CG and 96 IC flashes recorded by the Cordoba Marx Meter Array. The lightning data were in 3-D, 97 so they could employ the discharge height information to build a dataset with accurate 98 discharge types. They reported an overall accuracy of 97%. However, their proposed method 99 requires full waveform information, while most LLSs only retrieve a few parameters of 100 electric field waveforms of lightning discharges, making it somewhat difficult for existing 101 systems to adopt the method. 102

In this paper, we report a simple yet high-accuracy machine-learning technique based 103 on the Random Forest classifier to classify RSs. We will use a large dataset containing 104 about 15,000 negative RSs and many more IC discharges with accurate 3-D location information 105 to train and test the classifier. As will be described in this paper, many of the recorded 106 RSs and IC discharges produced atypical radiation waveforms that were challenging to 107 be correctly classified. However, the accuracy of our classifier is close to 99% demonstrated 108 by evaluations in various respects. Our classifier requires only some simple parameters 109 of lightning radiation waveforms, so it can be readily used by most LLSs. 110



Figure 1. (a) Negative RSs (black dots) observed from July 19 to August 26 in 2017. (b) Positive RSs observed from September 26, 2021 to September 3, 2022. Red squares represent observation sites of FALMA.

¹¹¹ 2 Observation and Data

During the summer of 2017, we set up a low-frequency (LF) lightning mapping system 112 called Fast Antenna Lightning Mapping Array (FALMA) in central Japan. The FALMA 113 consisted of 12 sites covering an area of about 80×80 km². Locations of these 12 sites are shown as red squares in Figure 1a. At every site, a fast antenna working in the frequency 115 band of 500 Hz to 500 kHz was used to receive radiation signals from lightning discharges. 116 The signals were recorded with a sampling rate of 25 MS/s. As described by Wu et al. 117 (2018a), thanks to improvements made in both the hardware and the software, we realized 118 high-quality 3-D lightning mapping with the FALMA. As can be seen from examples of 119 lightning flashes in Wu et al. (2018a) and Wu et al. (2019), 3-D mapping results of FALMA 120 have similar quality to those of very-high-frequency (VHF) systems such as the Lightning 121 Mapping Array (Rison et al., 1999).

Data obtained from July 19 to August 26 are used in this study for building and 123 testing the classifier for negative RSs. All data are reprocessed for this study. The largest 124 positive pulse (the same polarity as the negative RS, using the atmospheric electricity 125 sign convention) in each 20-ms window is located in 3-D. Only discharges located in the region shown in Figure 1a, a $90 \times 90 \text{ km}^2$ area over the FALMA network, are used in order 127 to ensure reliable 3-D locating. Pulses with source heights lower than 500 m are treated 128 as candidates of RSs. Their waveforms are then confirmed manually, and for some ambiguous 129 pulses, they are further manually located to determine their source heights. In this way, 130 we can unambiguously determine that the selected pulses are truly RSs. The number 131 of IC discharges are much larger than that of RSs, so we cannot manually confirm waveforms 132 of all IC discharges, and we only use pulses with source heights larger than 3000 m as 133 IC pulses. There are 14,898 pulses confirmed as negative RSs and 159,277 pulses as IC discharges. Locations of these RSs are shown as black dots in Figure 1a. It should be 135 noted that we will build a classifier for negative RSs rather than negative CG flashes; 136 a CG flash consists of at least one RS and also many IC discharges, both of which need 137 to be correctly classified. 138

Using the high-quality dataset of 2017 summer, we will establish the technique for 139 building the classifier as will be described in Sections 3.1 to 3.5. Further, using the same 140 technique, we will also build a classifier for positive RSs as will be described in Section 3.6. 141 However, positive RSs in central Japan in summer are quite rare (Wu et al., 2018b). In order to accumulate a large number of positive RSs, we will use the data collected during 143 a long period, from September 26, 2021 to September 3, 2022. During this period, we 144 set up a FALMA network covering a large area for 2-D locating of both summer and winter 145 lightning. Observation sites are shown as red squares in Figure 1b. A total of 8700 positive 146 RSs observed in an area with a radius of 300 km are identified and will be used for building 147 and testing the classifier for positive RSs. Locations of these positive RSs are shown as 148 black dots in Figure 1b. The procedure for the identification of these positive RSs will 149 be further described in Section 3.6. 150

Our classifiers will be built and tested mainly based on the Random Forest classifier,
 which is one of the most widely used machine-learning models for classification tasks.
 A brief comparison will also be made with the SVM classifier, another popular machine-learning model, in Section 3.4.

- ¹⁵⁵ 3 Methods and Results
- 156

3.1 Method to Evaluate the Performace of a Classifier

Before building the classifier, first we need to define some parameters as indicators of the performance of a classifier. One obvious parameter to evaluate the performance is the classification accuracy, or simply *accuracy*, that is, the percentage of true RSs in the waveforms classified as RSs. However, only this parameter is apparently not enough, as it is always possible to build a classifier with very strict criteria so that it only identifies very typical RS waveforms. Another important parameter is the identification efficiency, or simply *efficiency*, that is, the percentage of correctly classified RSs in all RSs.

Suppose the number of RSs is N_R , and the number of IC discharges is N_I . Of the N_R RSs, N_{Rc} are correctly classified (the subscript *c* stands for "correct"), and the remaining $N_R - N_{Rc}$ are misclassified as IC discharges. Of the N_I IC discharges, N_{Ic} are correctly classified, and the remaining $N_I - N_{Ic}$ are misclassified as RSs. The accuracy and the efficiency are defined as follows.

$$Accuracy = \frac{N_{Rc}}{N_{Rc} + (N_I - N_{Ic})} \tag{1}$$

169

$$Efficiency = \frac{N_{Rc}}{N_R} \tag{2}$$

During the process to build the classifier, we will experiment and tune various parameters of the classifier to make the accuracy and the efficiency as high as possible.

Normally a dataset is split into a larger training set and a smaller test set, with 172 the training set used to train a classifier and the test set used to test or evaluate the performance 173 of the classifier. In this study, we use an improved approach. All RS and IC data are combined, 174 shuffled and then divided into five equal parts. Each part is in turn used as the test set 175 and the remaining four parts combined are used as the training set. In this way, a classifier 176 is built and tested for five times and five results of accuracy and efficiency are calculated. The average values of five tests will be used as the final results. In this way, we can avoid 178 any random biases in the test set. Moreover, as will be described in Section 4, in this 179 way all data can be tested and we can find as many atypical waveforms as possible that 180 are difficult to be correctly classified. 181

¹⁸² 3.2 Waveform Paramaterization

We will define some waveform parameters to be used for building the classifier. First we describe the procedure to calculate waveform parameters based on multiple-site records. As waveforms recorded at a close distance contain the electrostatic and induction field components (e.g. Thottappillil et al., 1997) that may significantly distort the waveforms, observation sites within 40 km from a discharge are first excluded. Waveforms recorded by the remaining sites are used to calculate the parameters, and for each parameter, the median value of the results calculated based on these sites are used as the final result of the parameter for the discharge.

191

3.2.1 Parameters Related with Pulse Characteristics

First we define three basic parameters related with pulse characteristics. Definitions of these parameters are illustrated using an RS pulse in Figure 2a and an IC pulse in Figure 2b (blue parameters).

- 195 1. T_{rise} : The rise time of a pulse (10% to peak).
- **196** 2. T_{fall} : The fall time of a pulse (peak to zero).
- **3.** T_{half} : The pulse width at the half maximum.

With only these three basic parameters, we trained and tested the Random Forest 198 classifier using the negative RS and IC dataset obtained in 2017 summer. As described 199 in Section 3.1, the dataset is divided into five parts and each part in turn is used as the 200 test set, so the classifier is trained and tested for five times. The accuracy ranges from 201 72.25% to 73.57% with an average of 72.82%, and the efficiency ranges from 70.80% to 72.81% with an average of 71.59%. We also tried to add two related parameters, including 203 the pulse width, which is the sum of the rise time and fall time, and the ratio of fall time 204 to rise time, but the result has little difference (the average accuracy is 72.17% and the 205 average efficiency is 70.86%). 206

207 Indeed, with only these basic pulse parameters, it is difficult to accurately classify208 RSs.

209

3.2.2 Parameters Related with Relative Strength

An important feature of the RS waveform is that pulses right before and after an RS pulse is usually much weaker. The following parameters are defined to employ this feature. These parameters are also illustrated in Figures 2a and 2b.

- 1. R_{bp1} : The ratio of A_0 to A_{bp1} , in which A_0 is the peak amplitude of the target pulse, and A_{bp1} is the maximum amplitude of pulses right before the target pulse (from -100 μ s to 10% peak) as illustrated in Figure 2. The subscript *b* stands for "before", and the subscript *p* stands for "positive".
- 217 2. R_{bn1} , R_{bp2} , R_{an1} , R_{an1} , R_{ap2} , R_{an2} : These parameters are defined in the 218 same way as R_{bp1} , also illustrated in Figure 2. Note that the subscript *a* stands 219 for "after", and the subscript *n* stands for "negative".

The three parameters defined in Section 3.2.1 along with the eight new parameters defined above are used to train the Random Forest classifier. The accuracy of five tests ranges from 98.86% to 99.32% with an average of 99.02%, and the efficiency ranges from 98.02% to 98.66% with an average of 98.34%. It is clear that these new parameters representing the relative strength are very effective in the classification of RSs.



Figure 2. Illustration of waveform parameters using (a) an RS pulse and (b) an IC pulse. (c) Relative importance of waveform parameters.

3.2.3 Parameter Importance

The Random Forest classifier outputs a value indicating the relative importance 226 of each parameter in contributing to the performance, from which we can evaluate the 227 effectiveness of each parameter in the classification of RSs. The results are shown in Figure 2c. 228 Values of the importance of all parameters combined equal to 1. We can see that parameters 229 related with the pulse strength relative to previous pulses (red parameters in Figure 2) 230 are generally more important than other parameters. This is easy to understand as an 231 RS pulse is preceded by leader pulses which are usually much weaker than the RS pulse. 232 By contrary, an IC pulse is usually preceded by other IC pulses with comparable amplitudes. Therefore, parameters related with the relative strength are very effective in the classification 234 of RSs. 235

We can also see that parameters related with pulse characteristics (blue parameters) have relatively low importance, which is why the classifier performance is very poor with only these parameters as described in Section 3.2.1. It also indicates that traditional RS classification methods based on pulse characteristics are not very reliable.

3.3 Tradeoff Between Accuracy and Efficiency

From the above result, we can see one feature of the classifier is that the accuracy is always higher than the efficiency. It is obvious that increasing the efficiency usually implies decreasing the accuracy. However, it is desirable if we can control the tradeoff between the accuracy and the efficiency. For example, in some situations, it may be required to identify as many RSs as possible, so a high efficiency is essential while a low accuracy is tolerable. Next we will investigate two factors that influence the tradeoff between the accuracy and the efficiency.

248

240

3.3.1 Influence of Sample Size Imbalance

One reason for the higher accuracy in the classifier built in the previous section is a much larger sample of IC discharges compared with the sample of RSs. With such a biased dataset, the classifier is more likely to misclassify RSs, as also noted by Zhu et al. (2021). We can simply duplicate the sample of RSs to make the classifier identify more RSs, though at the cost of more misclassifications of IC discharges. Note that the duplication should only be made for the training set.

With the original dataset, 247 of 14,898 RSs (1.7%) are misclassified, but only 145 255 of 159,277 IC pulses (0.091%) are misclassified. If we duplicate the dataset of RSs in the 256 training set, the number of misclassified IC pulses increases to 162 while the number of 257 misclassified RSs decreases to 199. We tried to make more duplications and tested the 258 classifier, and the results of the accuracy and the efficiency are shown in Figure 3a. With 259 one duplication of the RS training set, the accuracy decreases from 99.02% to 98.91%but the efficiency increases from 98.34% to 98.66%. With two duplications, the accuracy decreases to 98.84% but the efficiency increases to 98.76%, very close to the accuracy. 262 With further duplications, we can see that both the accuracy and the efficiency are generally 263 very similar, changing between 98.75% and 98.85%, indicating that the sample size imbalance 264 does not have a significant effect any more. 265

If we use the average value of the accuracy and the efficiency as the indicator of the overall performance of a classifier, we can see from Figure 3a that with four duplications of the RS training set, the classifier has the highest performance with an accuracy of 98.84% and an efficiency of 98.81%. We treat this as the best performance of the classifier for negative RSs and this classifier will be used for further evaluations in the following section.



Figure 3. Variations of the accuracy and the efficiency with (a) different times of duplications of RS training set and (b) different thresholds of probability to classify RSs.

	Eleven Parameters (Section 3.2.2)		
Classifier	Accuracy (%)	Efficiency (%)	Time Cost (seconds)
Random Forest	99.02	98.34	20
SVM	98.43	97.42	96
Du	plicating RS Tra	ining Set (Section	n 3.3.1)
Classifier	Accuracy (%)	Efficiency (%)	Time Cost (seconds)
Random Forest	98.84	98.81	28
SVM	97.98	98.09	108

 Table 1. Comparison of the Random Forest classifier and the SVM classifier

271 3.3.2 Influence of Probability Thresholds

When classifying a pulse, the Random Forest classifier can output the probability
that the pulse is a true RS. By default, the classifier determines a pulse as an RS when
the probability is larger than 50%. By changing the probability threshold, we can conveniently
tune the accuracy-efficiency tradeoff.

Figure 3b shows variations of the accuracy and the efficiency related with the probability threshold. We can see that as the probability threshold increases, the accuracy increases while the efficiency decreases. This is easy to understand; a higher probability threshold represents stricter criteria to classify RSs, so naturally the identified RSs are more likely true RSs (higher accuracy), but at the same time fewer RSs can be identified (lower efficiency). In practice, when using the classifier we can set a customized probability threshold that fits the specific requirements of an application to achieve desired accuracy or efficiency.

283

3.4 Comparison of Different Machine-learning Models

Apart from the Random Forest classifier, another popular machine-learning model 284 for classification is the SVM classifier, which was used by Zhu et al. (2021) for the classification 70E of lightning pulses. Here we make a brief comparison of the Random Forest and the SVM classifiers. First we use the scheme described in Section 3.2.2 (using 11 parameters illustrated 287 in Figure 2) to train the classifiers, and the results are shown in Table 3.4 (upper part). 288 We can see the SVM classifier has slightly lower accuracy and efficiency than the Random 289 Forest classifier. Further, we use the scheme described in Section 3.3.1 (duplicating the 290 RS training dataset) to train the classifiers, and again, the SVM classifier has slightly 291 lower accuracy and efficiency. Another difference is in the time needed to train a classifier; 292 it takes less than 30 seconds to train an Random Forest classifier while the time needed 293 to train an SVM classifier is around 100 seconds. A significantly shorter time to build a classifier is potentially very useful as it would be more convenient to experiment various 295 combinations of parameters in order to boost the performance of the classifier. 296

297

3.5 Testing Using Remote Lightning Discharges

Lightning discharges used for training and evaluating classifiers described above are all very close to most of FALMA sites in order to ensure the 3-D location accuracy. However, many LLSs, especially national and continental LLSs, have long baselines of a few hundred kilometers, so lightning discharges observed by these systems are generally very far away from most observation sites. Therefore, it is desirable to evaluate the performance of a classifier for remote lightning discharges. We use lightning discharges located more than 150 km away from the center of the FALMA network in 2017 summer (the origin in Figure 1a) for this investigation. At such a large distance, only a small number of discharges can be located with sufficient accuracy, and we can only make 2-D locating, so we cannot classify RSs using the height information. Therefore, we manually inspected waveforms of all located events and determine their types.

There are a total of 594 located pulses. The classifier described in Section 3.3.1 (the 310 training set duplicated for four times) are used to classify these pulses. A total of 361 311 pulses were classified as RSs, and there was no clear misclassification. The remaining 233 312 pulses were classified as IC discharges, and four of them were likely RSs. However, it should 313 be noted that as there is no height information for these pulses, it is sometimes difficult 314 to determine the true discharge type, so it is possible that there were actually more misclassifications. 315 Assuming there are only four RSs misclassified as IC discharges, from Equations 1 and 316 2, we can get an accuracy of 100% and an efficiency of 98.9%. Note that when detecting 317 remote lightning discharges, as in the case of long-baseline LLSs, only a small portion 318 of IC discharges that are relatively strong can be located, so the chance of misclassifying an IC pulse as an RS is relatively low, which may be one reason for the 100% accuracy 320 in this evaluation. 321

The above results demonstrated that our classifier also has good performance when classifying remote RSs, so the classifier can also be used in long-baseline LLSs.

324 3.6 Classification of Positive Return Strokes

The methods described above can also be used to build a classifier for the classification 325 of positive RSs. However, positive CG flashes are very rare in summer thunderstorms 326 in central Japan. As reported by Wu et al. (2018b), only 46 positive CG flashes consisting 327 of 53 positive RSs were observed and could be located in 3-D during the summer observation 328 of 2017. Therefore, here we also include the data obtained in other periods. First we use the 690 positive RSs observed during the winter of 2018 (Wu et al., 2022) to build a preliminary 330 classifier for the identification of positive RSs. Then we use this classifier to search the 331 data recorded in about one year from September of 2021 for possible positive RSs. As 332 described in Section 2, during this period, we set up a FALMA network with long baselines 333 for 2-D locating of both summer and winter lightning. Waveforms of the identified positive 334 RSs by the preliminary classifier are manually confirmed to exclude obvious false classifications. 335 Indeed, the preliminary classifier identified many pulses that were clearly IC pulses and 336 we painstakingly excluded all apparent IC pulses by manual inspections. In this way, we 337 collected the data of 8700 positive RSs, locations of which are shown in Figure 1b. Note 338 that there is no height information for these positive RSs, so this dataset is not as accurate 339 as the negative RS dataset in 2017 summer used in previous sections. 340

For IC data, we also use the data of summer observation of 2017 as these data have 341 accurate 3-D location results. However, different from the IC dataset for the negative 342 RS classifier, IC pulses for building positive RS classifier should have the same polarity 343 as positive RSs. So we located IC pulses having the same polarity as positive RSs and 344 selected those with heights larger than 3 km, the same treatment as that in building the 345 negative RS classifier. On the other hand, as the size of positive RS dataset is relatively 346 small, we do not need too many IC data, so for simplicity, we only located one IC pulse 347 in every 50-ms window. Finally, we collected a total of 113,922 IC pulses. 348

Using these datasets, and with the same scheme for building negative RS classifier described in Section 3.3.1, we built and tested the classifier for positive RSs. It is found that with the RS training set duplicated for one time, the classifier has the best overall performance. It has an accuracy of 99.04% and an efficiency of 98.37%, generally similar to the performance of the negative RS classifier. This result demonstrated that as long as there are enough data of positive RSs, we can also build a high-accuracy classifier for positive RSs in the same way as building the negative RS classifier.

Although the dataset of positive RSs does not have 3-D location information and thus is not as accurate as the negative RS dataset, as positive RSs are much rarer than negative RSs and it is very difficult to collect a large and reliable sample, we believe our classifier is very valuable for future observations and researches. Moreover, as all waveforms of identified positive RSs have been manually confirmed, the classifier likely has an accuracy similar to that of the manual classification.

³⁶² 4 Atypical Intracloud and Return Stroke Waveforms

As described in Section 3.1, the whole dataset of 2017 summer is divided into five parts, with each part in turn used as the testing set and the remaining four parts combined used as the training set. In this way, all pulses can be tested and we can identify as many pulses as possible that are potentially difficult to classify. Using the classifier built in Section 3.3.1 with the RS training set duplicated for four times, all pulses are classified. Of the 14,898 RSs, 178 were misclassified as IC discharges, and of the 159,277 IC pulses, 173 were misclassified as RSs. Waveform figures of all these misclassified pulses are provided in the data repository.

There are several common reasons for misclassifications of RS pulses as IC pulses. 370 Waveforms of four examples are shown in Figures 4a-d, all of which are misclassified as 371 IC discharges and whose source heights have been confirmed to be close to the ground. 372 First, it is well known that RSs striking tall grounded objects usually produce very narrow 373 pulses (Araki et al., 2018; Cai et al., 2022; Pavanello et al., 2007; Zhu et al., 2018), making 374 it easy to misclassify them as IC discharges. One example is shown in Figure 4a. This pulse is located near a transmission tower, and its pulse width is only about 6 μ s, indicating 376 that it is likely produced by an RS striking the tower. Second, two RSs sometimes occur 377 sequentially with a very small time difference of a few tens of microseconds, and if the 378 second RS has a larger peak than the first one, the second RS may be misclassified as 379 an IC pulse. One example is shown in Figure 4b. Such RSs are likely the so-called "multiple-termination 380 strokes" (Kong et al., 2009; Sun et al., 2016) or "forked strokes" (Ballarotti et al., 2005), 381 with two RSs induced by two branches of the same leader. Third, an RS may occur almost simultaneously with IC discharges of other lightning flashes, resulting in a peculiar waveform 383 and thus misclassified as an IC discharge. One example is shown in Figure 4c. While the 384 positive pulse is confirmed to be produced by an RS, the two negative pulses labeled as 385 "IC" are produced by IC discharges in an independent lightning flash and are located about 386 87 km away from the RS. The resultant waveform appears to be abnormal and is difficult 387 to be identified as an RS. Finally, some RSs apparently produce waveforms that are largely 388 different from typical RS waveforms but the reason is not yet clear. One example is shown in Figure 4d. The pulse has a rise time of about 18 μ s while its fall time is only about 9 μs . 391

Another example of an RS producing abnormal waveform is shown in Figure 5 along with location results of the preceding leader. This RS is a subsequent RS. We can see from the location results in Figure 5a a dart leader with a speed of about 4×10^6 m/s preceding the RS, and the RS is located very close to the ground as indicated by the cross sign. From Figure 5c, we can see details of the RS waveform. It contains two peaks with the second peak much larger than the first one, resulting in a much larger rise time than the fall time. Without the 3-D location results, it is very difficult to determine that the waveform is produced by an RS.

The major reason for IC discharges misclassified as RSs is that waveforms of some
IC discharges have some similar features as those of RSs. Four examples are shown in
Figures 4e-h. All of these waveforms appear very similar to those of RSs. However, their
source heights range from 5.9 to 15.1 km, indicating that they are produced by IC discharges.



Figure 4. (a)-(d) Atypical E-change waveforms produced by RSs but misclassified as IC discharges. (e)-(h) Atypical E-change waveforms produced by IC discharges but misclassified as RSs. The value of d represents the distance between the discharge and the observation site recording the waveform. The value of h represents the source height of the IC discharge.



Figure 5. Location result and E-change waveforms of an RS misclassified as an IC discharge. (a) Height-time location results of the dart leader preceding the RS. The cross sign represents the RS. (b) E-change waveform of the RS and preceding discharges. (c) E-change waveform of the RS. The value of d represents the distance between the RS and the observation site recording the waveform.

We also manually located these pulses to make sure that there were no large errors in the source height results. We can see that these pulses have relatively short rise times and much longer fall times. Pulses in Figures 4e-g also have fine structures superimposed on the falling part, similar to waveforms of first RSs, and the pulse in Figure 4h resembles the waveform of a subsequent RS. These similar features as RS waveforms make it almost impossible to correctly classify them as IC discharges without the 3-D location results.

Another example of an IC pulse appearing similar to RS pulses is shown in Figure 6 410 along with location results of preceding discharges. From the height-time location results in Figure 6a, we can see that a leader first propagated above 6.5 km and then descended 412 to a height of about 5.5 km, and then the large IC pulse is produced, represented by the 413 cross sign. From the E-change waveform in Figure 6c, we can see that the large IC pulse 414 is very similar to an RS pulse, with preceding pulses resembling stepped leader pulses. 415 With the help of the 3-D location results, we can be sure that this RS-like pulse is produced 416 by IC discharges. We are not aware of any study reporting such RS-like IC pulses. In 417 our future studies, we will explore the mechanism responsible for these special IC pulses. 418

These examples of special RS and IC waveforms illustrate the fact that some RSs and IC discharges produce atypical radiation waveforms from which the discharge types cannot be accurately determined, resulting in a fundamental difficulty to achieve very high accuracy and efficiency using only waveform information. This result also illustrates the importance of accurate 3-D location results in scientific investigations of lightning phenomena.

425 5 Conclusions

Using a large dataset with 3-D location results, we built a classifier for radiation 426 waveforms of negative RSs based on the Random Forest classifier. Eleven simple parameters 427 are defined for building the classifier, including three parameters related with pulse characteristics 128 and eight parameters related with relative strength of pulses. A classification accuracy 429 of 98.84% and an identification efficiency of 98.81% are achieved. We also demonstrated 430 methods to tune the tradeoff between the accuracy and the efficiency so the classifier can 431 be used in applications with different requirements of the accuracy or the efficiency. Although 432 the classifier is built based on the observation of a compact lightning mapping system, 433 we demonstrated that the classifier also has high accuracy and efficiency for remote lightning 434 discharges and can be readily used in long-baseline LLSs. With the same methods, we 435 also built a classifier for positive RSs which has similarly high accuracy and efficiency as the classifier for negative RSs. 437

Misclassified RS and IC waveforms are examined and some common reasons for misclassifications
are analyzed. We demonstrated that RSs sometimes produce radiation waveforms that
are largely different from normal RS waveforms, and IC discharges sometimes produce
waveforms that appear very similar to RS waveforms. Therefore, some RS and IC waveforms
are fundamentally difficult to be correctly classified without 3-D location information,
and it is likely that such misclassifications commonly exist in most LLSs. The results
also imply the importance of 3-D location results in detailed analyses of lightning phenomena.

445 Open Research Section

Datasets for building and testing the classifiers as well as waveform figures of all positive and negative RSs can be found at https://doi.org/10.5281/zenodo.7641792. Sample Python scripts for using the classifiers will be made publicly available after the acceptance of this paper.



Figure 6. Location result and E-change waveforms of an IC pulse misclassified as an RS pulse. (a) Height-time location results of the IC pulse and preceding discharges. The cross sign represents the location of the IC pulse. (b) E-change waveform of the IC pulse and preceding discharges. (c) E-change waveform of the IC pulse. The value of *d* represents the distance between the IC discharge and the observation site recording the waveform.

450 Acknowledgments

This study was supported by the Ministry of Education, Culture, Sports, Science, and Technology of Japan (Grants 20H02129 and 21K03681).

453 References

454	Araki, S., Nasu, Y., Baba, Y., Rakov, V. A., Saito, M., & Miki, T. (2018, sep).
455	3-d finite difference time domain simulation of lightning strikes to the 634-m
456	tokyo skytree. Geophysical Research Letters, $45(17)$, $9267-9274$. doi:
457	$\frac{1}{10000000000000000000000000000000000$
458	observations of negative ground flashes on a millisecond-scale <i>Ceonhysical</i>
459	Research Lettere 32(23) doi: https://doi.org/10.1020/2005gl023880
460	Betz H D Schmidt K Oettinger P $\&$ Wirz M (2004 jun) Lightning
401	detection with 3-d discrimination of intracloud and cloud-to-ground discharges
402	Geophysical Research Letters 31(11) n/a-n/a doi: 10.1029/2004gl019821
464	Cai, L., Liu, W., Zhou, M., Wang, J., Yan, R., Tian, R., & Fan, Y. (2022, dec).
465	Differences of electric field parameters for lightning strikes on tall towers and
466	nonelevated objects. <i>IEEE Transactions on Electromagnetic Compatibility</i> ,
467	64(6), 2113–2121. doi: https://doi.org/10.1109/temc.2022.3207237
468	Cummins, K. L., & Murphy, M. J. (2009, aug). An overview of lightning locating
469	systems: History, techniques, and data uses, with an in-depth look at the u.s.
470	NLDN. IEEE Transactions on Electromagnetic Compatibility, 51(3), 499–518.
471	doi: https://doi.org/10.1109/temc.2009.2023450
472	Fleenor, S. A., Biagi, C. J., Cummins, K. L., Krider, E. P., & Shao, XM.
473	(2009, feb). Characteristics of cloud-to-ground lightning in warm-season
474	thunderstorms in the central great plains. Atmospheric Research, 91(2-4),
475	333–352. doi: https://doi.org/10.1016/j.atmosres.2008.08.011
476	Kohlmann, H., Schulz, W., & Pedeboy, S. (2017, oct). Evaluation of EUCLID
477	IC/CG classification performance based on ground-truth data. In 2017 intermetional summasium on lightning metaction (VIV CIDDA) IEEE doi:
478	https://doi.org/10.1100/sindo.2017.8116806
479	Kong X Oio X Zhao V & Zhang T (2000 feb) Characteristics of
480	negative lightning flashes presenting multiple-ground terminations on a
482	millisecond-scale. Atmospheric Research, 91(2-4), 381–386. doi: https://
483	doi.org/10.1016/j.atmosres.2008.03.025
484	Leal, A. F., Rakov, V. A., & Rocha, B. R. (2019, aug). Compact intracloud
485	discharges: New classification of field waveforms and identification by lightning
486	locating systems. Electric Power Systems Research, 173, 251–262. doi:
487	https://doi.org/10.1016/j.epsr.2019.04.016
488	Lin, Y. T., Uman, M. A., Tiller, J. A., Brantley, R. D., Beasley, W. H., Krider,
489	E. P., & Weidman, C. D. (1979). Characterization of lightning return stroke
490	electric and magnetic fields from simultaneous two-station measurements.
491	Journal of Geophysical Research, 84 (C10), 6307. doi: https://doi.org/10.1029/
492	jc084ic10p06307
493	Paul, C., Heidler, F. H., & Schulz, W. (2020, feb). Performance of the european
494	lightning detection network EUCLID in case of various types of current
495	purses from upward fighting measured at the persenderg tower. IEEE Transactions on Floatromagnetic Commetibility $60(1)$ 116 122
496	$\frac{1100-123}{100} \text{ doi } 00000000000000000000000000000000000$
497	Pavanello D. Rachidi F. Janischewsky W. Rubinstein M. Hussein A. M.
498	Petrache E Jaquier A (2007 jul) On return stroke currents and remote
₩ ₽ ₽	electromagnetic fields associated with lightning strikes to tall structures.
501	2. experiment and model validation. Journal of Geophysical Research:
502	Atmospheres, 112(D13). doi: https://doi.org/10.1029/2006id007959

503	Rison, W., Thomas, R. J., Krehbiel, P. R., Hamlin, T., & Harlin, J. (1999, dec). A
504	GPS-based three-dimensional lightning mapping system: Initial observations
505	in central new mexico. Geophysical Research Letters, 26(23), 3573–3576. doi:
506	https://doi.org/10.1029/1999gl010856
507	Schulz, W., Diendorfer, G., Pedeboy, S., & Poelman, D. R. (2016, mar). The
508	european lightning location system EUCLID – part 1: Performance analysis
509	and validation. Natural Hazards and Earth System Sciences, 16(2), 595–605.
510	doi: https://doi.org/10.5194/nhess-16-595-2016
511	Sun, Z., Qie, X., Liu, M., Jiang, R., Wang, Z., & Zhang, H. (2016, jan).
512	Characteristics of a negative lightning with multiple-ground terminations
513	observed by a VHF lightning location system. Journal of Geophysical
514	Research: Atmospheres, 121(1), 413–426. doi: https://doi.org/10.1002/
515	2015id023702
516	Thottappillil, R., Rakov, V. A., & Uman, M. A. (1997, mar). Distribution of charge
517	along the lightning channel: Relation to remote electric and magnetic fields
517	and to return-stroke models <i>Journal of Geophysical Research</i> : Atmospheres
510	102(D6) 6987–7006 doi: https://doi.org/10.1029/96id03344
519	Wang I Huang O Ma O Chang S He I Wang H Gao C (2020 feb)
520	Classification of VLF/LF lightning signals using sensors and deen learning
521	mathedes Sensore $20(A)$ 1030 doi: https://doi.org/10.3300/s200/1030
522	Wu T Wang D Huang H l_r Takagi N (2021 nov) The strongest negative
523	lightning strokes in winter thunderstorms in inpan
524	Lettere /8(21) doi: https://doi.org/10.1020/2021gl005525
525	$E_{\rm Letters}$, 40 (21). doi: https://doi.org/10.1029/2021gi095525
526	fast antonnas — <i>Coonhusical Research Letters</i> 15(8) 3608-3705 — doi: https://
527	doi org/10.1002/2018gl077628
528	$W_{\rm H}$ T Wang D fr Takagi N (2018b aug) L coating proliminary breakdown
529	pulses in positive cloud to ground lightning <i>Lowrnal of Coonducting Research</i> :
530	Atmospheres, doi: https://doi.org/10.1020/2018id028716
531	$W_{\rm H}$ T Wang D & Takagi N (2010 sop) Velocities of positive leaders
532	in intracloud and nogative cloud to ground lightning flashes
533	of Coonhusical Research: Atmospheres 19/(17.18) 0083-0005
534	https://doi.org/10.1020/2010jd030783
535	Wu T Wang D & Takagi N (2021 aug) Compact lightning strokes in winter
536	thunderstorms Lowrnal of Coonhusical Research: Atmospheres 196(15) doi:
537	https://doi.org/10.1020/2021id034032
538	Wu T Wang D k Takagi N (2022 nov) On the intensity of first return
539	strokes in positive cloud-to-ground lightning in winter Lowrnal of Geophysical
540	Research: Atmospheres 197(22) doi: https://doi.org/10.1020/2022id037282
541	Zhu V Bitzer P Bakov V & Ding Z (2021 jan) A machine-learning approach
542	to classify cloud-to-ground and intracloud lightning Geophysical Research
543	Letters $\frac{1}{2}$ doi: https://doi.org/10.1020/2020gl001148
544	Zhu V Bakov V A Tran M D Lyu W $\&$ Micu D D (2018 sep) A
545	modeling study of parrow electric field signatures produced by lightning strikes
540	to tall towars Journal of Geophysical Research: Atmospheres 199(18) doi:
547	https://doi.org/10.1020/2018id028016
548	The V Bakey V A Tran M D & Nag A (2016 dec) A study of national
549	lightning detection network responses to netural lightning based on ground
550	truth data acquired at LOC with emphasis on cloud discharge activity
551	Journal of Geophysical Research: Atmospheres 191(94) 14 651-14 660
552	bttps://doi.org/10.1002/2016jd025574
553	Zhu V Stock M Lapierre I & DiCangi E (2022 may) Ungrades of the earth
554	networks total lightning network in 2021 Remote Sensing 1/(0) 2200 doi:
555	https://doi org/10/2300/rs1/002200
550	noops.//uonorg/10.0000/1814002200