



*universe*

IMPACT  
FACTOR  
**2.5**

CITESCORE  
**4.3**

Article

---

# Convolutional Neural Network Processing of Radio Emission for Nuclear Composition Classification of Ultra-High-Energy Cosmic Rays

---

Tudor Alexandru Calafeteanu, Paula Gina Isar and Emil Ioan Slușanschi

Special Issue

Advanced Studies in Ultra-High-Energy Cosmic Rays

Edited by

Dr. Gina Isar and Prof. Dr. François Montanet



<https://doi.org/10.3390/universe10080327>

## Article

# Convolutional Neural Network Processing of Radio Emission for Nuclear Composition Classification of Ultra-High-Energy Cosmic Rays

Tudor Alexandru Calafeteanu <sup>1,2</sup> , Paula Gina Isar <sup>1,\*</sup>  and Emil Ioan Slușanschi <sup>2</sup> 

<sup>1</sup> Institute of Space Science—Subsidiary of INFLPR, 077125 Bucharest-Magurele, Romania; tudor.calafeteanu@stud.acs.upb.ro

<sup>2</sup> Faculty of Automatic Control and Computer Science, National University of Science and Technology Politehnica Bucharest, 060042 Bucharest, Romania

\* Correspondence: gina.isar@spacescience.ro

**Abstract:** Ultra-high-energy cosmic rays (UHECRs) are extremely rare energetic particles of ordinary matter in the Universe, traveling astronomical distances before reaching the Earth's atmosphere. When primary cosmic rays interact with atmospheric nuclei, cascading extensive air showers (EASs) of secondary elementary particles are developed. Radio detectors have proven to be a reliable method for reconstructing the properties of EASs, such as the shower's axis, its energy, and its maximum ( $X_{\max}$ ). This aids in understanding fundamental astrophysical phenomena, like active galactic nuclei and gamma-ray bursts. Concurrently, data science has become indispensable in UHECR research. By applying statistical, computational, and deep learning methods to both real-world and simulated radio data, researchers can extract insights and make predictions. We introduce a convolutional neural network (CNN) architecture designed to classify simulated air shower events as either being generated by protons or by iron nuclei. The classification achieved a stable test error of 10%, with Accuracy and  $F_1$  scores of 0.9 and an MCC of 0.8. These metrics indicate strong prediction capability for UHECR's nuclear composition, based on data that can be gathered by detectors at the world's largest cosmic rays experiment on Earth, the Pierre Auger Observatory, which includes radio antennas, water Cherenkov detectors, and fluorescence telescopes.

**Keywords:** cosmic rays; air showers; radio detection; deep learning



**Citation:** Calafeteanu, T.A.; Isar, P.G.; Slușanschi, E.I. Convolutional Neural Network Processing of Radio Emission for Nuclear Composition Classification of Ultra-High-Energy Cosmic Rays. *Universe* **2024**, *10*, 327. <https://doi.org/10.3390/universe10080327>

Academic Editor: Lorenzo Iorio

Received: 18 July 2024

Revised: 9 August 2024

Accepted: 11 August 2024

Published: 15 August 2024



**Copyright:** © 2024 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

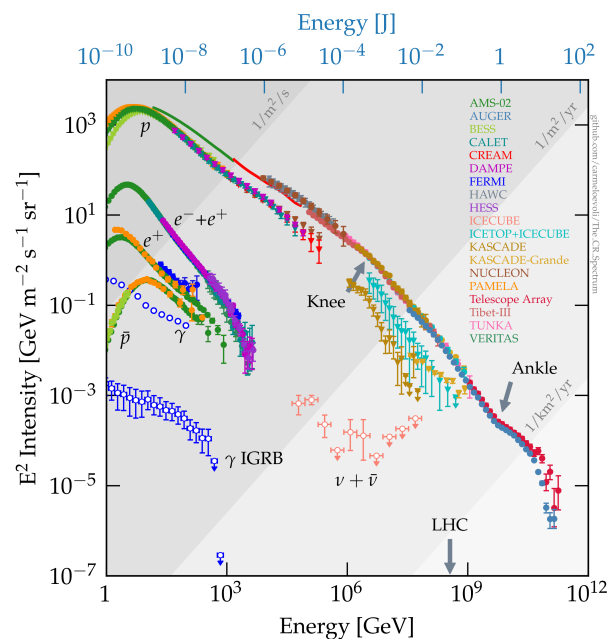
## 1. Introduction

Cosmic rays [1,2] are high-energy charged particles that originate from outer space. They are composed of light nuclei (such as proton and helium nuclei, at energies around  $10^{18}$  eV (1 EeV)) up to heavy nuclei (such as iron, at increasing energies) [3]. Their energies [1] can reach up to  $10^{20}$  eV ( $\approx 16$  Joules), larger than the energy levels reached in the most advanced particle accelerator on Earth, the Large Hadron Collider (LHC) at CERN, where the maximum collision energy attainable is 14 TeV [4]. Various astrophysical processes accelerate these particles, enabling them to reach such high energies. Based on their origins, cosmic rays are classified as [5] solar energetic particles (originating from the Sun, mainly solar flares), galactic cosmic rays (originating from our galaxy), or extra-galactic cosmic rays (originating from external galaxies). Possible galactic and extra-galactic sources are accretion shocks in large-scale structures, active galactic nuclei, gamma-ray bursts, and neutron stars or magnetars [6].

Cosmic rays are important agents in multi-messenger astronomy—the coordinated observation and interpretation of signals from various information carriers, called messengers, such as electromagnetic radiation, gravitational waves, neutrinos, and cosmic rays [5,7–10]. For example, during solar flares, both electromagnetic radiation and cosmic

rays are emitted [11]. Their coordinated signals reveal important information that aids researchers in identifying their source.

As the energy increases, the flux of cosmic rays decreases rapidly (Figure 1), where the flux of particles with energies above  $10^{20}$  eV is less than one particle per  $\text{km}^2$  per century. Two important transition features are observed: the knee ( $\approx 3 \cdot 10^{15}$  eV) and the ankle ( $\approx 10^{19}$  eV). The knee feature is related to galactic cosmic rays, while the latter marks the shift from galactic to extra-galactic origins. Due to their low flux, only a few events have been detected for cosmic rays with energies around 100 EeV. The highest-energy cosmic ray ever observed was detected in 1991, with an estimated energy of 320 EeV. The flux of cosmic rays is strongly suppressed above the cutoff ( $\approx 5 \cdot 10^{19}$  eV), but there is no consensus on exactly how much this limit is intrinsic to their sources and how much it is due to energy losses during the intergalactic propagation [12,13].



**Figure 1.** The flux of cosmic rays with data from various experiments [14]. Reproduced with permission from Carmelo Evoli. Published by Zenodo, 2020.

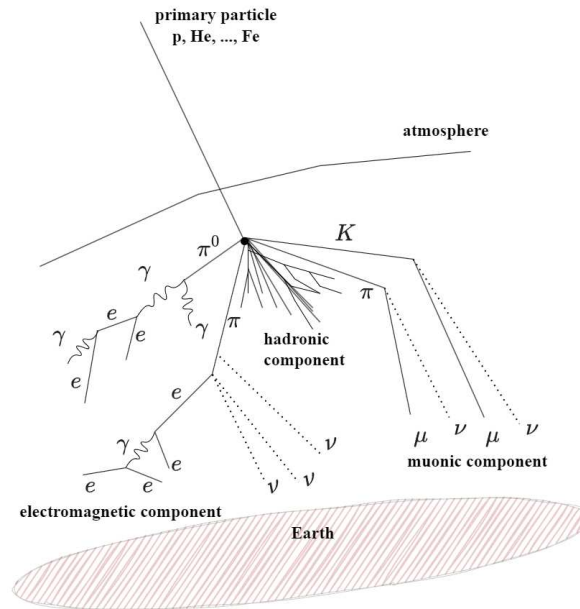
Direct measurements of cosmic rays with energies over  $10^{14}$  eV are challenging [1]. This is due to both the decreasing flux with increasing energy and the fact that ultra-high-energy cosmic rays (UHECRs), which are particles with energies exceeding 1 EeV, can pass through detectors placed at the top of the atmosphere without undergoing interactions.

Indirect measurements are preferred for UHECRs. They are conducted by measuring not the primary particle itself but the secondary particles produced in an extensive cascade, following the initial collision of the cosmic ray with air molecules in the atmosphere. This process is known as an *extensive air shower* (EAS) (Figure 2), which consists of hadronic, muonic, and electromagnetic components. The electromagnetic component, composed of electrons, positrons, and photons, carries about 90% of the primary particle’s energy [15], making it the most valuable component for energy measurements. The remaining 10% is carried by muons and neutrinos.

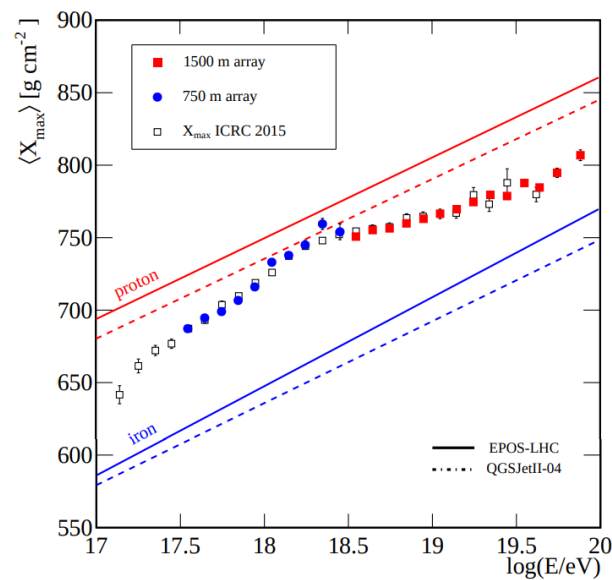
After the initial collision of a cosmic ray with air molecules in our atmosphere, secondary particles begin to be produced in a cascade manner. Their number increases until it reaches a maximum at a certain atmospheric depth, called  $X_{\text{max}}$ . Beyond this depth, the number of secondary particles starts decreasing, due to ionization losses [15].

The footprint of the EAS provides valuable information about the cosmic ray, including its energy, arrival direction, core position on the ground, and the shower maximum ( $X_{\text{max}}$ ). The reconstruction of  $X_{\text{max}}$ , i.e., the depth in the atmosphere at which the energy deposition

in the shower is greatest, is essential for determining the mass compositions of the cosmic rays (Figure 3). These measurements can be compared with  $X_{\max}$  prediction from air shower simulations for the two extremes of primary particle types (proton and iron nuclei) for different hadronic interaction models.



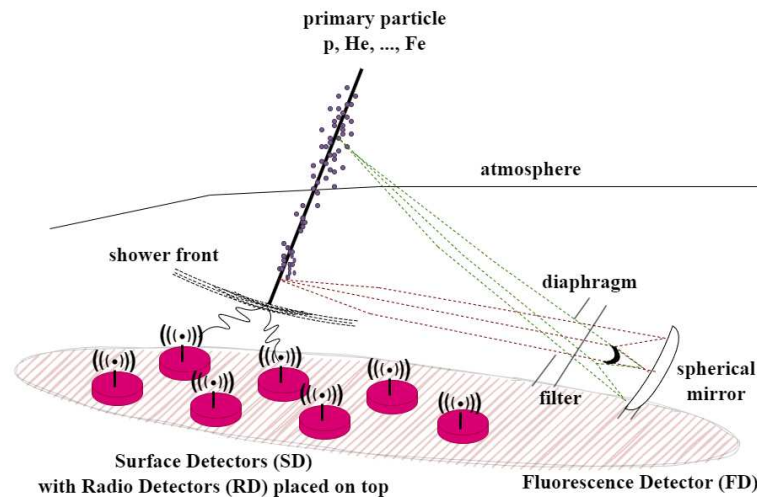
**Figure 2.** Illustration of a developing air shower with its main components: electromagnetic, muonic, and hadronic. Inspired by and adapted from [16].



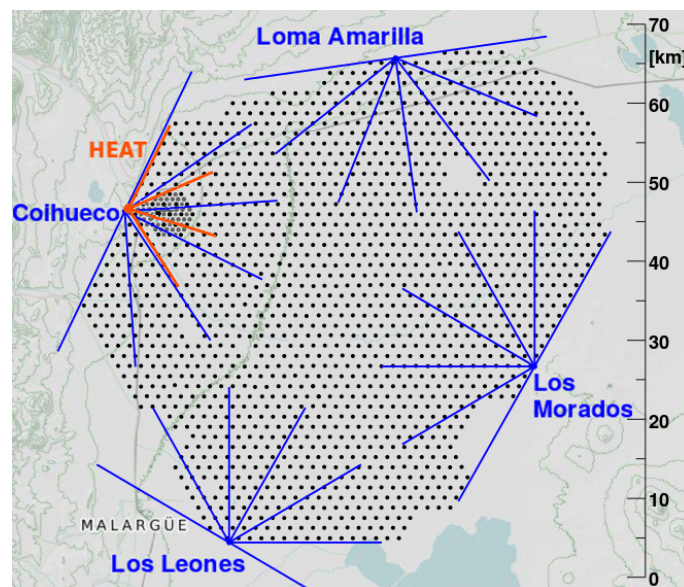
**Figure 3.** Mean values of  $X_{\max}$  as function of energy, for data measured by the fluorescence and surface detectors of the Pierre Auger Observatory, in comparison with hadronic models [17].

The Pierre Auger Observatory [18] is the world’s largest cosmic ray experiment, operating in the western Mendoza province, Argentina since 2004. Its main goal is to understand the origin and nature of UHECRs [19]. The observatory is unique, in that it combines complementary detection techniques (Figure 4); 27 fluorescence telescopes capture the ultraviolet light emitted by excited nitrogen molecules in the atmosphere; 1660 water Cherenkov detectors, which are ground-based tanks filled with purified water overlooked by three photo-multipliers that detect Cherenkov radiation, are dispersed over

an area of 3000 km<sup>2</sup> with a 1.5 km spacing between them; additionally, a radio antenna array is sensitive in the frequency range of 30–80 MHz. Starting from an initial array of only 153 radio antennas (the Auger Engineering Radio Array), the Auger upgrade (AugerPrime) [20] will include, among other extensions, a radio antenna on top of each ground detector, enabling electric field detection on the full map (Figure 5).



**Figure 4.** Illustration of the hybrid detection of an air shower with the surface and the fluorescence detectors, as implemented in the Pierre Auger Observatory. Inspired by and adapted from [21].



**Figure 5.** Real map of the Pierre Auger Observatory, where the red and blue lines indicate the field of view of the fluorescence telescopes and the black points represent the surface particle detectors [18].

As the flux of UHECRs decreases with increasing energy, Monte Carlo cosmic ray air shower simulation software, such as CORSIKA 7 (COsmic Ray SIMulations for KASCADE) [22], are used to generate a statistically significant number of datasets. Various radio antennas response studies have been conducted on CoREAS [23] (CORSIKA-based Radio Emission from Air Showers) simulations, part of CORSIKA, with antennas from the full Auger map, in the order of tens or even a few hundreds [24,25].

The main emission mechanisms that cause the electric field recorded by the antennas are the geomagnetic and Askaryan effects [26]. The former represents the dominant contribution, and it depends on the geomagnetic position of the experiment, while the

latter is caused by the negative charge excess developed in the shower, and it depends on the shower direction.

The radio pulses captured by antennas provide valuable data for reconstructing shower observables [27], such as energy,  $X_{\max}$ , and the zenith and azimuth angles [25,28]. This is because the radio footprint left by an EAS depends significantly on the properties of the UHECR. For example, inclined showers with a larger zenith angle produce an elongated, elliptical footprint, whereas vertical showers with a small zenith angle generate a more circular footprint. The azimuth angle determines the rotation of the footprint's geometry. Additionally, the energy of the UHECRs plays a crucial role in the energy fluence and electric field's strength, captured by radio antennas [29].

Machine learning (ML) has become widely used nowadays in physics. Depending on their complexity, machine learning models are classified as shallow and deep. Shallow learning algorithms, such as linear regression, logistic regression, decision trees, and support vector machines, are more suited to simpler tasks that require fewer parameters and less computational power. They often excel in scenarios where the relationship between input features and output is relatively straightforward and the dataset is not overly large.

On the other hand, deep learning algorithms leverage neural networks with many layers, which enables them to model complex, non-linear relationships in data. These deep neural networks are particularly well-suited for tasks that involve high-dimensional data and intricate patterns, such as image recognition, natural language processing, and speech recognition. Deep learning models, such as convolutional neural networks (CNNs), require substantial computational resources and large datasets to perform effectively, but at the same time they offer significant improvements in Accuracy and capability for these complex tasks.

Since the radio footprint is highly correlated with many properties of UHECRs, it can be effectively used in image recognition techniques to infer the nature of the primary particle. This is where deep learning, particularly CNNs, plays an important role. By leveraging convolutional layers that use information from neighboring pixels, CNNs can detect patterns in images and emphasize the energy deposit and the geometry of the radio footprint left by an induced EAS.

There have been several studies on deep learning processing of radio data using CNNs—such as energy estimation, using ZHAireS simulations [30], and classification between signal and background noise [31,32], using CoREAS simulations.

We propose a convolutional neural network architecture that classifies simulated CoREAS events between those generated by protons and those by iron nuclei. The model was trained and evaluated on a dataset of  $\approx 3000$  events ( $\approx 2000$  with proton and  $\approx 1000$  with iron) from the Pierre Auger Collaboration for this purpose. The paper is arranged as follows: in Section 2, we describe the datasets used in the analysis, the radio imaging techniques employed, and the distribution of the numerical features. Section 3 presents a first shallow machine learning analysis on the numerical features only and the performance of the algorithms employed. Section 4 describes the architecture of the neural network, the training process, and the metrics used to assess the model's performance. We conclude in Section 5 with an evaluation of the model's success in the classification task.

## 2. Data Exploration and Preprocessing

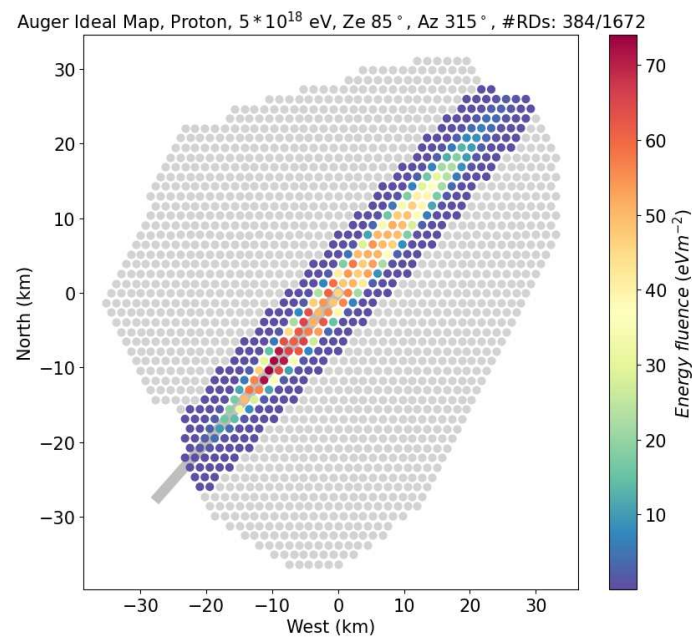
Different signal properties are used in radio analysis, such as the signal-to-noise ratio (SNR), which indicates how effectively the signal can be distinguished from noise; the electric field strength, measured in  $\mu\text{V}/\text{m}$ , which describes the voltage induced by the signal into a one-meter-long radio antenna; and the energy fluence, measured in  $\text{eV}\text{m}^{-2}$ , which describes the energy deposit over the unit area. These properties can be analyzed for each polarization component of the electric field, or as a magnitude (for the SNR and the electric field strength), or as a sum (for energy fluence) of all components.

We considered first our dataset of a small simulation library done with CORSIKA v7.7420 and the CoREAS option for two types of primary particles (proton and iron) with an

energy of  $5 \cdot 10^{18}$  eV, for different air shower geometries defined by zenith (from  $70^\circ$  to  $85^\circ$ , in steps of  $5^\circ$ ) and azimuth (from  $0^\circ$  to  $360^\circ$ , in steps of  $45^\circ$ ) angles. This resulted in a total of 32 simulated events per primary particle. The simulated radio pulses were conducted for a fixed shower core in the center of the array and for all radio antennas located within a certain distance of the shower axis (in the shower plane),  $r_{\max} < \max(4 \cdot r_0(\text{zenith}), 1500 \text{ m})$ , where  $r_0$  is the radius of the Cherenkov ring [33]. The high-energy hadronic interaction model was QGSJETII-04; the low-energy interaction model was UrQMD, with a thinning parameter of  $1 \times 10^{-6}$ . The atmospheric condition was representative of the Pierre Auger Observatory in October at Malargüe, with the corresponding magnetic field at the experiment location.

The CoREAS simulations in our dataset were performed on different subsets of radio antennas from the complete ideal Auger map of  $3000 \text{ km}^2$ , at  $1400 \text{ m a.s.l.}$  This ideal map filled the gaps from the real Auger map, ensuring that all the radio antennas were evenly spaced. In Auger Phase II, each point also represents a radio antenna. The time series of the simulated electric field, captured by each radio antenna, was further processed, to calculate the energy fluence (Equation (1) [34], Figure 6):

$$F = \frac{1}{Z_0} \int_{-T}^T E(t)^2 dt, \quad (1)$$



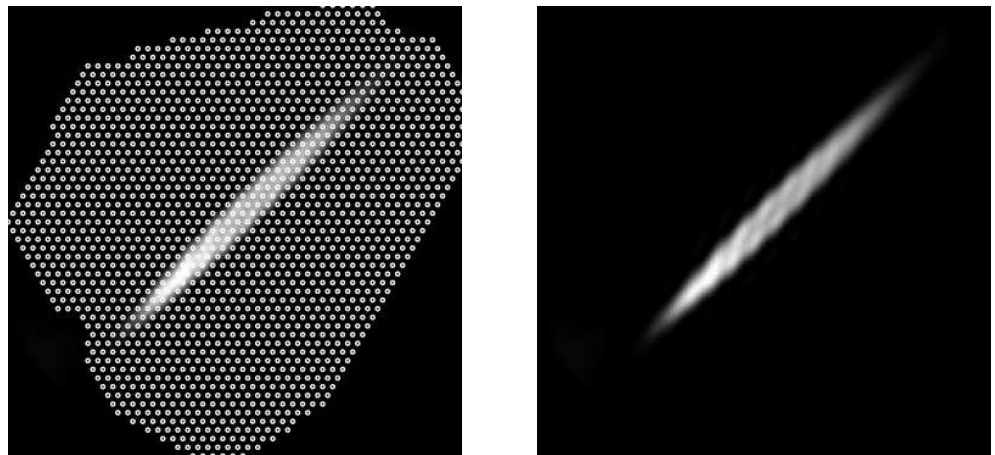
**Figure 6.** The energy fluence for an ideal Auger map simulation for a proton UHECR with an energy of  $5 \cdot 10^{18}$  eV, a zenith of  $85^\circ$ , and an azimuth of  $315^\circ$  (an EAS coming from the south-east). The radio antennas that were not simulated are depicted in light gray. The shower axis is depicted in gray.

The plot can be further transformed in a gray-scale image, by applying a linear interpolation function to the energy fluence on the whole map viewed as a mesh grid of  $400 \times 400$  points (Figure 7, left). The radio antennas that were not simulated were considered to have zero energy fluence over that area. After removing the radio antennas used for the interpolation from the plot, only the energy fluence footprint is shown (Figure 7, right).

The main issue with this imaging technique is that the energy fluence—as well as other signal properties, like maximum amplitude—leaves a relatively small footprint on the full map. This is because the full Auger map is best suited for UHECRs that are more inclined. For zenith angles less than  $80^\circ$ , the radio footprint becomes very narrow, almost point-like (Figure 8).

We employed four types of imaging techniques:

- Max local method: Each radio antenna's energy fluence is MinMax scaled, where the maximum value is determined per simulation. This method presents the radio footprint on each unit area, relative to the footprint of the entire simulation. This method is presented in Figures 7 and 8.
- Log max local method: Similar to the Max local method, but with a  $\log_{10}$  transformation applied to the energy fluence, to enhance the visibility of the energy deposit comparison between different area units.
- Max global method: Each radio antenna's energy fluence is MinMax scaled, where the maximum value is determined across all simulations, which is over  $4.24 \cdot 10^5 \text{eVm}^{-2}$ , captured by an antenna from a simulation of a vertical iron UHECR with an energy over  $10^{20}$  eV, coming from the south-east. This method presents the radio footprint on each unit area, relative to the footprint of the entire dataset, aiding in the comparison of energy deposits between simulations from the whole dataset.
- Log max global method: Similar to the Max global method, but with a  $\log_{10}$  transformation applied to the energy fluence, to enhance the visibility of the energy deposit comparison between different area units.



**Figure 7.** The interpolation of energy fluence over the ideal Auger map for simulation depicted in Figure 6 with (left) and without (right) radio antennas depicted.



**Figure 8.** The energy fluence footprint for an ideal Auger map simulation for a proton UHECR with an energy of  $5 \cdot 10^{18}$  eV and an azimuth of  $315^\circ$  (left—zenith  $80^\circ$ , right—zenith  $75^\circ$ ).

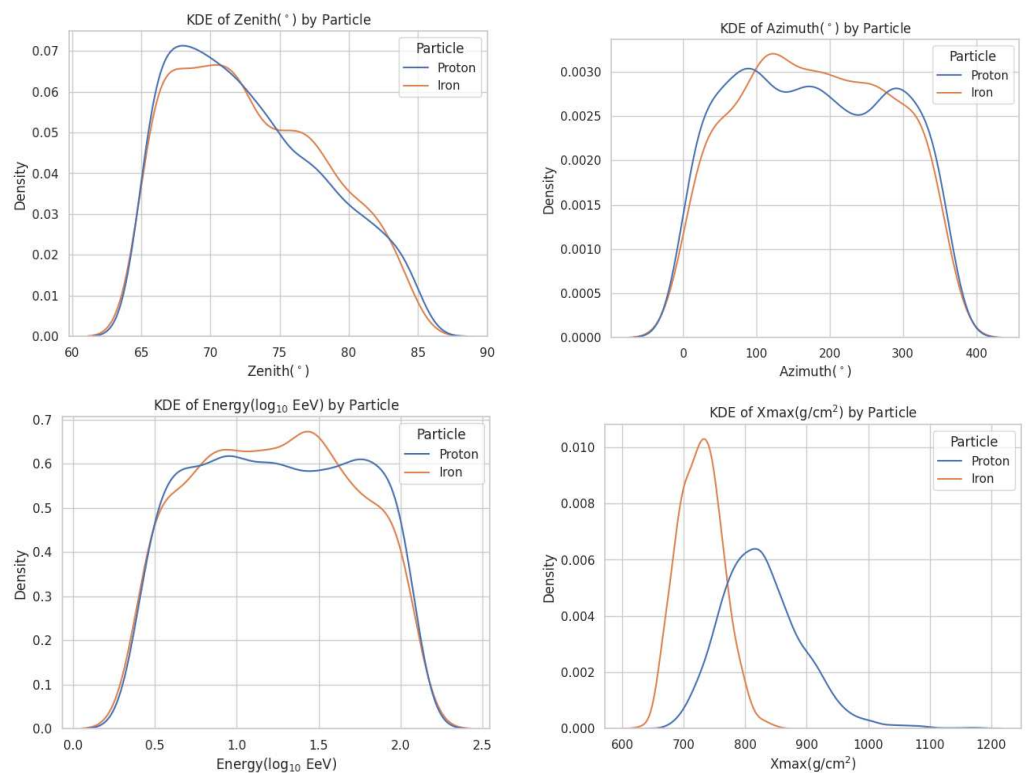
In order to increase our dataset from the order of tens to at least thousands, we used air showers simulation libraries from the Pierre Auger Collaboration [33] (Auger

simulation dataset). These simulations were done with CORSIKA/CoREAS v7.7401 for two primary particles: approximately 2000 showers for protons and about 1000 showers for iron. The core positions were randomly distributed uniformly over the entire ideal Auger layout. The simulations included zenith angles between  $65^\circ$  and  $85^\circ$ , azimuth angles between  $0^\circ$  and  $360^\circ$ , and primary particle energies in  $\log_{10}(E/eV)$  from 18.4 to 20.1. The high-energy hadronic interaction model was Sibyll-2.3d; the low-energy interaction model was UrQMD. The thinning and atmospheric conditions, as well as the antennas selection algorithm, were the same as those used in our dataset. We further applied the same four radio imaging techniques.

Apart from the aforementioned images, we extracted four observables that significantly impacted the radio footprint: the zenith, azimuth, energy, and  $X_{\max}$  (Table 1). Figure 9 shows a nearly uniform distribution for the first three, while  $X_{\max}$  followed a more Gaussian distribution. In order to normalize the independent features, we applied the MinMax method to the former three and the Z-score to the latter.

**Table 1.** Summary statistics for zenith, azimuth, energy, and  $X_{\max}$ .

Observable	Min	Max	Mean
Zenith ( $^\circ$ )	65.0	84.99	72.94
Azimuth ( $^\circ$ )	0.12	359.94	178.54
Energy (EeV)	2.51	119.72	30.01
$X_{\max}$ (g/cm $^2$ )	628.92	1174.7	791.42



**Figure 9.** Kernel density estimation (KDE) for (top left) zenith, (top right) azimuth, (bottom left) energy, and (bottom right)  $X_{\max}$ .

The Pearson correlation coefficient (PCC) also highlighted the importance of  $X_{\max}$  in determining the nuclear composition of the UHECRs (Figure 10). While the other observables had a low impact, due to their uniform distribution for both primaries, there was a stronger linear correlation ( $-0.64$ ) between  $X_{\max}$  and the primary particle type (proton encoded as 0, and iron as 1): proton-induced showers reached their maximum at a greater depth in the atmosphere, compared to those induced by iron. Additionally,

a significant linear correlation (0.36) existed between  $X_{\max}$  and energy: air showers induced by cosmic rays with higher energies reached their maximum at greater depths in the atmosphere.

The Auger simulation dataset was split into training and testing sets using a 70/30 ratio, with stratification by the particle feature (proton or iron), to ensure balanced representation.

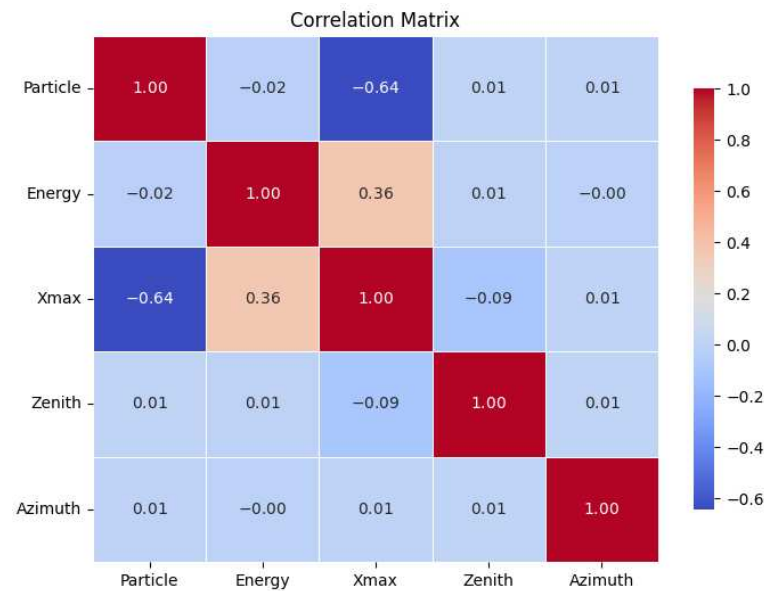


Figure 10. Linear correlation of the features (proton encoded as 0, and iron as 1).

### 3. Shallow Learning Classification

To understand the impact of radio imaging techniques on a deep learning model, we first evaluated how the numerical data performed in the classification task, using several shallow classification algorithms from the *scikit-learn* [35] Python library. The models considered in our study included decision tree (*DecisionTreeClassifier*), random forest (*RandomForestClassifier*), support vector machine—SVM (*SVC*), k-nearest neighbors (*KNeighborsClassifier*), logistic regression (*LogisticRegression*), gradient boosting (*GradientBoostingClassifier*), and Gaussian naive Bayes (*GaussianNB*). Each model was initialized with a random state of 42, to ensure reproducibility. The hyper-parameters for each model were tuned using grid search, and the best parameters were selected based on performance metrics.

The evaluation metrics used to provide a comprehensive assessment of the classification performance were Accuracy (Equation (2)), which highlights the ratio between correct and total predictions, but can be misleading in class-imbalanced datasets [36]; MCC (Matthews correlation coefficient, Equation (3)), which is more reliable than Accuracy for class-imbalanced datasets [36]; and  $F_1$  score (Equation (4)), which balances precision and recall without considering true negatives. While the Accuracy and  $F_1$  scores range from 0 to 1 (with higher values indicating better predictions), MCC ranges from  $-1$  to 1, where  $-1$  indicates total disagreement between prediction and actual class, 0 indicates randomness, and 1 indicates a perfect match. Given our class-imbalanced dataset, we used the MCC as the primary metric to assess both the hyper-parameter tuning of the model and the comparison between different models.

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN} \tag{2}$$

$$\text{MCC} = \frac{TP \cdot TN - FP \cdot FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}} \tag{3}$$

$$F_1 = \frac{2 \cdot TP}{2 \cdot TP + FP + FN} \quad (4)$$

**Table 2.** Best three models and their optimized hyper-parameters for classification between proton and iron nuclei using different features as input. The model with the best MCC score is bolded.

Input Features	Model	Best Parameters	Accuracy	F1 Score	MCC
$X_{\max}$	<b>logistic regression</b>	<b>C: 0.1,</b> <b>penalty: l1,</b> <b>solver: liblinear</b>	<b>0.8590</b>	<b>0.7988</b>	<b>0.6904</b>
	gradient boosting	learning_rate: 0.01, max_depth: 3, n_estimators: 100, subsample: 0.8	0.8568	0.7925	0.6841
	SVM	C: 10, gamma: auto, kernel: rbf	0.8525	0.7982	0.6828
$X_{\max},$ Energy	<b>SVM</b>	<b>C: 10,</b> <b>gamma: scale,</b> <b>kernel: linear</b>	<b>0.9214</b>	<b>0.8902</b>	<b>0.8291</b>
	logistic regression	C: 10, penalty: l1, solver: liblinear	0.9203	0.8886	0.8266
	gradient boosting	learning_rate: 0.1, max_depth: 3, n_estimators: 100, subsample: 0.8	0.9139	0.8758	0.8105
$X_{\max},$ Energy, Zenith, Azimuth	<b>logistic regression</b>	<b>penalty: l2,</b> <b>solver: saga</b>	<b>0.9300</b>	<b>0.9014</b>	<b>0.8472</b>
	random forest	criterion: entropy, max_depth: 10, min_samples_leaf: 1, min_samples_split: 10, n_estimators: 100	0.9193	0.8841	0.8226
	gradient boosting	learning_rate: 0.1, max_depth: 3, n_estimators: 50, subsample: 0.8	0.9182	0.8816	0.8199

Table 2 summarizes, based on their MCC score, the best three models with their optimized hyper-parameters and the resulting performance metrics, including MCC,  $F_1$ -score, and Accuracy for different features as input. The logistic regression and the support vector machine with linear kernel performed the best in the classification task, indicating that the data had a clear decision boundary that could be effectively captured by these linear classifiers to predict the nuclear composition of the UHECRs. The importance of  $X_{\max}$  in the classification was evidenced by the MCC score of 0.69, which denoted good performance. Adding the energy feature increased the MCC to 0.83, while also including the spherical angles zenith and azimuth gave a slight boost in performance, up to about 0.85.

However, in reality, the values of the  $X_{\max}$ , energy, zenith, and azimuth angles are not known and need to be reconstructed from the detector response data. Since  $X_{\max}$  is an indispensable value in nuclear composition prediction and has been shown to be reliably reconstructed using deep neural network processing of data from the water Cherenkov detectors from the Pierre Auger Observatory [37], we investigated the possibility of replacing the energy and the zenith and azimuth angles with the energy fluence of the electric field captured by radio antennas, as the radio footprint is significantly influenced by these

observables (as explained in Section 1). The proposed deep learning model is a CNN that receives as input the radio imaging techniques discussed in Section 2, along with the  $X_{\max}$ , and classifies simulated air shower events between those induced by protons and those by iron nuclei.

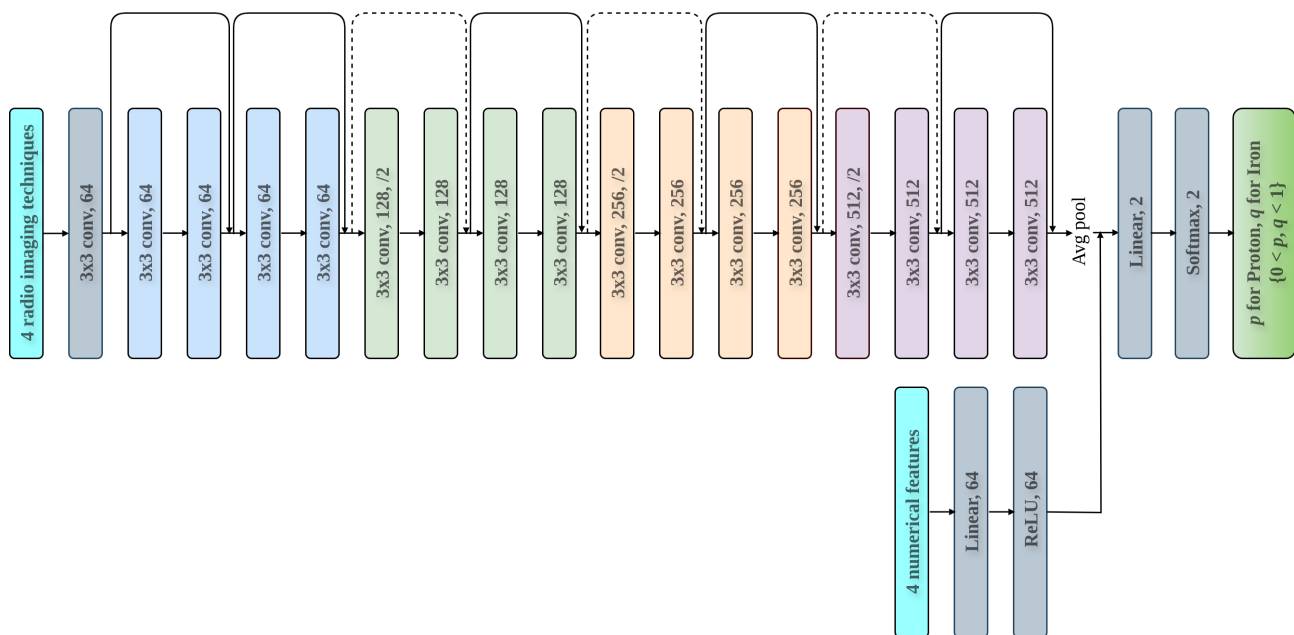
#### 4. Deep Learning Classification

We further discuss the architecture of the neural network, its training process, and, finally, the evaluation of its performance.

The CNN is a slightly modified version of the pretrained ResNet-18 [38]. As the ResNet-18 expects three-channel RGB images as input, we replace its first layer with a convolutional layer that receives the four radio images. The rest of the parameters, i.e., the number of output channels, the kernel size, the stride, the padding, and the bias are preserved. Its last layer (the classification layer) is also replaced with the identity operator, to integrate particle-specific features, such as zenith, azimuth, energy, and  $X_{\max}$ .

The numerical features (up to four) are linearly transformed into a higher-dimensional space (the default number of dimensions is 64). The rectified linear unit (ReLU) activation function is then applied, to introduce non-linearity, enabling more complex learning.

The images and numerical features processing produce two output vectors, which are further concatenated and linearly transformed into a final vector. This vector is converted using the softmax function into a probability distribution of two outcomes, resulting in a  $p$  probability that the UHECR particle is a proton, and a  $q$  probability that it is an iron nucleus. The final architecture of the CNN model is displayed in Figure 11.



**Figure 11.** Architecture of the CNN. The convolutional layers and average pooling operation are part of the original ResNet 18 architecture. Inspired by and adapted from [39].

The dataset used for training included only the  $X_{\max}$  feature and the four different techniques of energy fluence imaging. To artificially increase the training dataset, the data augmentation technique *random resized crop* was used. It randomly selects a portion of an image, resizes it, and then crops it to a specified size of  $224 \times 224$  pixels. The selected portion has a scale between 0.75 and 1, and it has an aspect ratio between 0.75 and 1.33, providing variability in the training data.

We used the *CrossEntropyLoss* for the loss function and the *AdamW* optimizer with a learning rate of 0.001 and a weight decay of 0.05. The model parameters were updated every mini-batch of 64 images. Since the training set had 2166 images, this resulted in a

total of  $\lceil \frac{2166}{64} \rceil = 34$  gradient descent steps per epoch. An epoch is one complete pass of the training data through the algorithm. The training data were reshuffled at every epoch.

At each epoch, we plotted the training and test errors, the loss curve, and the MCC,  $F_1$  and Accuracy scores, for a total of 50 epochs. The test errors (Equation (5)) for both proton and iron (Figure 12) followed a similar decreasing trend, stabilizing around 10%. This demonstrates the robust ability of the model to predict the nuclear composition of the UHECR based on only  $X_{\max}$  observable and radio images.

$$\text{Error rate} = \frac{FP + FN}{TP + TN + FP + FN} = 1 - \text{Accuracy} \tag{5}$$

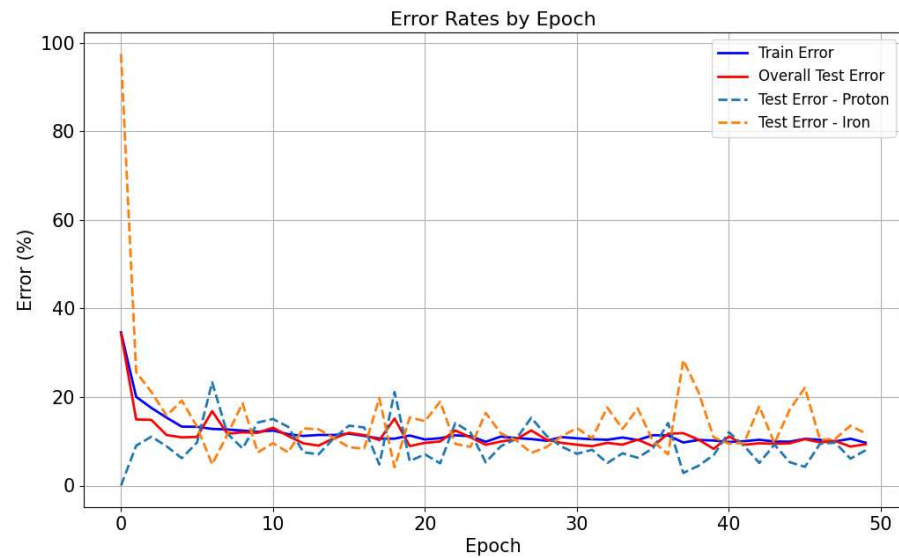


Figure 12. The error rates by epoch.

The training and testing loss (Figure 13) both decreased steadily, indicating that the model was learning effectively, stabilizing to about 0.25.

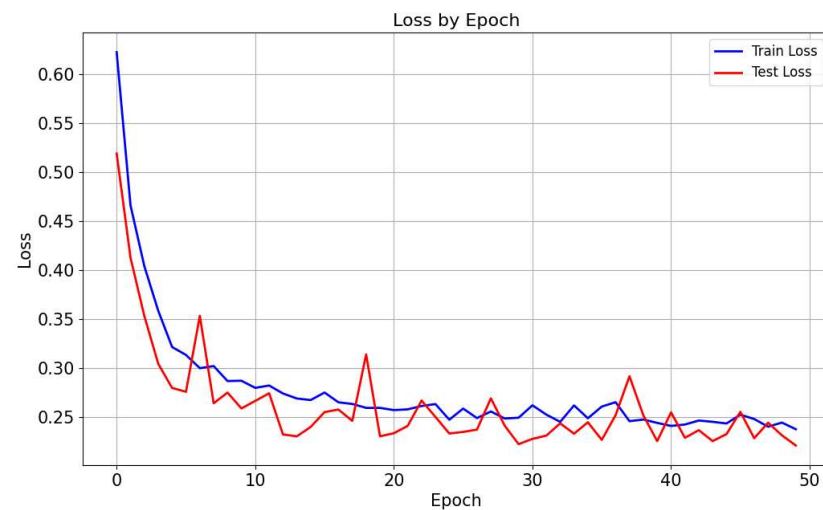
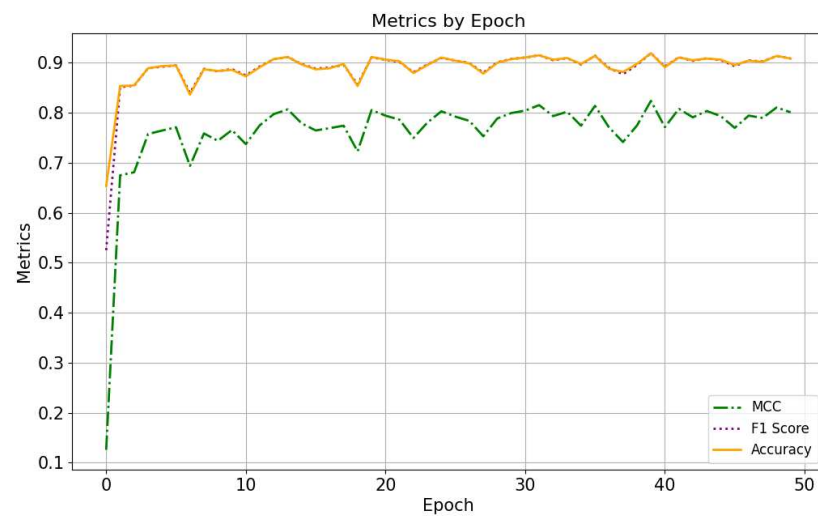


Figure 13. The loss rates by epoch.

Regarding the evaluation metrics (Figure 14), the MCC increased rapidly in the first few epochs and then gradually stabilized around 0.8. This indicated a strong correlation between predicted and actual values. Accuracy increased quickly and then stabilized around 0.9, suggesting good performance on the classification task. The lower MCC value

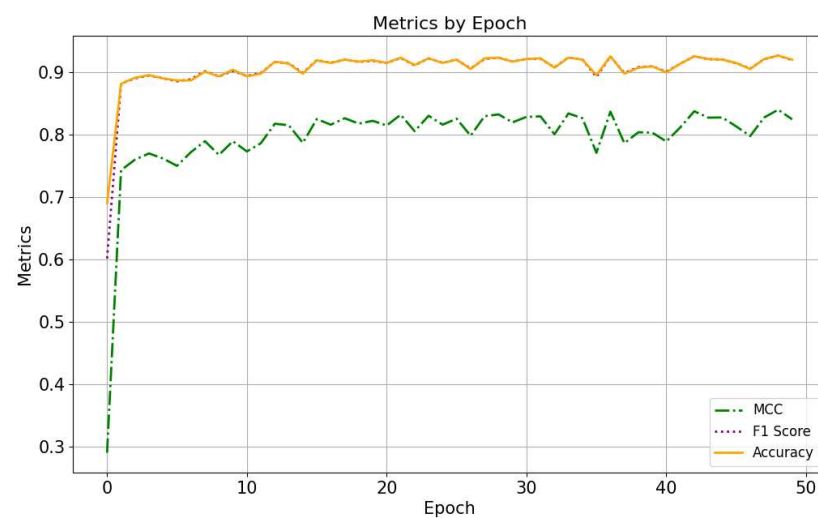
compared to Accuracy may indicate that the class imbalance affected the predictions. The  $F_1$  score increased sharply at the beginning and stabilized around 0.9, indicating a good balance between precision and recall.



**Figure 14.** The evaluation metrics by epoch.

Since the logistic regression model using all four numerical features as input achieved an MCC of about 0.85 (as shown in Table 2) and the CNN, which used the  $X_{\max}$  and energy fluence images as input, converged to an MCC of 0.8 (as seen in Figure 14), it is evident that the neural network model demonstrates good performance. This indicates that its classification capability is comparable to a scenario where all the other observables (energy and zenith and azimuth angles) are already known.

We ran the training again, this time including all four numerical features ( $X_{\max}$ , energy, zenith angle, and azimuth angle). The evaluation metrics (Figure 15) showed a slight overall performance increase: the MCC stabilized just above 0.8, while the Accuracy and  $F_1$  scores rose just above 0.9. This shows that the radio images already contained information about these newly included features, illustrating that the radio imaging techniques effectively substituted for the unknown energy of the primary particle and the air shower direction.



**Figure 15.** The evaluation metrics by epoch for all four numerical features included ( $X_{\max}$ , energy, zenith angle, and azimuth angle).

## 5. Conclusions

We successfully implemented a convolutional neural network model to classify UHE-CRs as either iron nuclei or protons, based on the simulated radio data from their induced air showers, using the ideal map of the Pierre Auger Observatory. The model used the  $X_{\max}$  observable, which has a strong linear relationship with the primary particle type, along with various radio imaging techniques to fully leverage the convolutional layers of the model. It showed a stable test error of 10%, with Accuracy and  $F_1$  scores of 0.9 and an MCC of 0.8. Comparing to an MCC of 0.85, when all the other desired observables—i.e., the energy and the spherical angles (zenith and azimuth)—are known and used by a logistic regression model, this indicates that the energy fluence images were successfully able to replace the unknown observables that influenced the radio footprint in the nuclear composition classification task. Moreover, the information about these unknown observables was strongly contained in the radio images, since additionally including them as inputs in the neural network model did not significantly increase the model's performance.

As future work, we aim to improve the model performance and reduce the test error to below the current 10% by employing additional methods, such as new radio imaging techniques and more hyper-parameter tuning. Enhancements, like adding more layers or using three outputs for the energy fluence (one per polarization channel, instead of the current sum), could be implemented. The dataset could also be increased by including more types of cosmic rays, and an even more interesting approach would be to use reconstructed radio data from the real experiment, to test the model's Accuracy.

**Author Contributions:** Conceptualization, T.A.C., P.G.I. and E.I.S.; methodology, T.A.C.; software, T.A.C.; validation, T.A.C. and P.G.I.; formal analysis, T.A.C.; investigation, T.A.C. and P.G.I.; resources, P.G.I.; data curation, T.A.C.; writing—original draft preparation, T.A.C.; writing—review and editing, P.G.I. and E.I.S.; visualization, T.A.C. and P.G.I.; supervision, P.G.I. and E.I.S.; project administration, P.G.I.; funding acquisition, P.G.I. All authors have read and agreed to the published version of the manuscript.

**Funding:** This work was supported by the Romanian Ministry of Research, Innovation, and Digitization, CNCS-UEFISCDI, project number PN-III-P1-1.1-TE-2021-0924/TE57/2022, within PNCDI III, and under the Romanian National Core Program LAPLAS VII-contract no. 30N/2023.

**Data Availability Statement:** Parts of the datasets analyzed in the current study are owned by the Pierre Auger Collaboration. Approval from the corresponding author is required for access upon reasonable request.

**Acknowledgments:** Tudor Alexandru Calafeteanu would like to express sincere gratitude to his colleague and friend, David-Gabriel Ion, for his valuable feedback and advice throughout this project. Additionally, Paula Gina Isar would like to thank the Pierre Auger Collaboration for the simulation library used in the machine learning training (Sections 3 and 4) and for useful feedback on the manuscript.

**Conflicts of Interest:** The authors declare no conflicts of interest.

## References

1. Glaser, J.C. Absolute Energy Calibration of the Pierre Auger Observatory Using Radio Emission of Extensive Air Showers. Ph.D. Thesis, RWTH Aachen University, Aachen, Germany, 2017.
2. Isar, P.G.; Hirnea, D.; Jipa, A. Cosmic rays air showers properties and characteristics of the emitted radio signals using analytical approaches and full Monte Carlo simulations. *Rom. Rep. Phys.* **2020**, *72*, 301.
3. Hoerandel, J.R. [The Pierre Auger Collaboration]. The nature and origin of ultra high-energy cosmic rays. *Europhys. News* **2012**, *43*, 24–27. [[CrossRef](#)]
4. Evans, L.; Bryant, P. LHC machine. *J. Instrum.* **2008**, *3*, S08001. [[CrossRef](#)]
5. Adamo, M.; Pietroni, S.; Spurio, M. Astrophysical sources and acceleration mechanisms. *arXiv* **2022**, arXiv:2202.09170.
6. Kotera, K.; Olinto, A.V. The astrophysics of ultrahigh-energy cosmic rays. *Annu. Rev. Astron. Astrophys.* **2011**, *49*, 119–153. [[CrossRef](#)]
7. Kampert, K.-H. et al. [The Pierre Auger Collaboration] Multi-Messenger Physics with the Pierre Auger Observatory. *Front. Astron. Space Sci.* **2019**, *6*, 00024.

8. Batista, R.A.; Biteau, J.; Bustamante, M.; Dolag, K.; Engel, R.; Fang, K.; Kampert, K.-H.; Kostunin, D.; Mostafa, M.A.; Murase, K.; et al. Open Questions in Cosmic Ray Research at Ultrahigh Energies. *Front. Astron. Space Sci.* **2019**, *6*, 23.
9. Meszaros, P.; Fox, D.B.; Hanna, C.; Murase, K. Multi-messenger astrophysics. *Nat. Rev. Phys.* **2019**, *1*, 585–599. [[CrossRef](#)]
10. Coleman, A.; Eser, J.; Mayotte, E.; Sarazin, F.; Schröder, F.G.; Soldin, D.; Venters, T.M.; Aloisio, R.; Alvarez-Muñiz, J.; Alves Batista, R.; et al. Ultra high energy cosmic rays. The intersection of the Cosmic and Energy Frontiers. *Astropart. Phys.* **2023**, *149*, 102819.
11. Bazilevskaya, G.A. Once again about origin of the solar cosmic rays. *J. Phys. Conf. Ser.* **2017**, *798*, 012034. [[CrossRef](#)]
12. Aab, A. et al. [The Pierre Auger Collaboration] Features of the energy spectrum of cosmic rays above  $2.5 \times 10^{18}$  eV using the Pierre Auger Observatory. *Phys. Rev. Lett.* **2020**, *125*, 121106. [[CrossRef](#)]
13. Aab, A. et al. [The Pierre Auger Collaboration] Measurement of the cosmic-ray energy spectrum above  $2.5 \times 10^{18}$  eV using the Pierre Auger Observatory. *Phys. Rev. D* **2020**, *102*, 062005. [[CrossRef](#)]
14. Evoli, C. The Cosmic-Ray Energy Spectrum. *Zenodo* **2020**. [[CrossRef](#)]
15. Anchordoqui, L.A. Ultra-high-energy cosmic rays. *Phys. Rep.* **2019**, *801*, 1–93. [[CrossRef](#)]
16. Travnicek, P. Detection of High-Energy Muons in Cosmic Ray Showers. Ph.D. Thesis, Charles University, Prague, Czech Republic, 2004.
17. Aab, A. et al. [The Pierre Auger Collaboration]. Inferences on mass composition and tests of hadronic interactions from 0.3 to 100 EeV using the water-Cherenkov detectors of the Pierre Auger Observatory. *Phys. Rev. D* **2017**, *96*, 122003. [[CrossRef](#)]
18. Gora, D. [The Pierre Auger Collaboration]. The Pierre Auger Observatory: Review of latest results and perspectives. *Universe* **2018**, *4*, 128. [[CrossRef](#)]
19. The Pierre Auger Collaboration. The Pierre Auger Cosmic Ray Observatory. *NIM-A* **2015**, *798*, 172–213. [[CrossRef](#)]
20. Pierre Auger Collaboration. The Pierre Auger Observatory and its upgrade. *Sci.-Rev.-End World* **2020**, *1*, 8–33. [[CrossRef](#)]
21. Obermeier, A. The Fluorescence Yield of Air Excited by Electrons Measured with the AIRFLY Experiment. Ph.D. Thesis, FZKA 7284, Karlsruhe University, Karlsruhe, Germany, 2007.
22. Heck, D.; Knapp, J.; Capdevielle, J.N.; Schatz, G.; Thouw, T. *CORSIKA: A Monte Carlo Code to Simulate Extensive Air Showers*; FZKA Report 6019; Forschungszentrum Karlsruhe: Karlsruhe, Germany, 1998.
23. Huege, T.; Ludwig, M.; James, C.W. Simulating radio emission from air showers with CoREAS. *AIP Conf. Proc.* **2013**, *1535*, 128–131.
24. Isar, P.G.; Hirnea, D. The response of a model hexagonal detector area to radio signals from ultra-high energy cosmic rays air showers. *Rom. Rep. Phys.* **2022**, *74*, 301.
25. Abdul, H.A. et al. [The Pierre Auger Collaboration] Radio measurements of the depth of air-shower maximum at the Pierre Auger Observatory. *Phys. Rev. D* **2024**, *109*, 022002. [[CrossRef](#)]
26. Isar, P.G. Radio signals from highly energetic extensive air showers: Status and new prospective. *Rom. Rep. Phys.* **2023**, *75*, 301.
27. Aab, A. et al. [The Pierre Auger Collaboration] Observation of inclined EeV air showers with the radio detector of the Pierre Auger Observatory. *JCAP* **2018**, *10*, 026.
28. Aab, A. et al. [The Pierre Auger Collaboration] Energy estimation of cosmic rays with the Engineering Radio Array of the Pierre Auger Observatory. *Phys. Rev. D* **2016**, *93*, 122005. [[CrossRef](#)]
29. Gaté, F. [The Pierre Auger Collaboration]. Radio detection of cosmic rays with the Auger Engineering Radio Array. In Proceedings of the 25th European Cosmic Ray Symposium, Turin, Italy, 4–9 September 2016.
30. Beatriz de Souza Pancrácio de Errico. Deep Learning-Based Energy Reconstruction of Cosmic Rays with Radio Emission Simulations. Master's Thesis, Institute of Physics-UFRJ, Rio de Janeiro, Brazil, 2023.
31. Erdmann, M.; Schlüter, F.; Šmída, R. Classification and recovery of radio signals from cosmic ray induced air showers with deep learning. *J. Instrum.* **2019**, *14*, P04005. [[CrossRef](#)]
32. Rehman, A.; Coleman, A.; Schröder, F.G.; Kullgren, D.; Abbasi, R.; Ackermann, M.; Adams, J.; Agarwalla, S.; Aguilar, J.; Ahlers, M.; et al. Search for Cosmic-Ray Events Using Radio Signals and CNNs in Data from the IceTop Enhancement Prototype Station. In Proceedings of the 38th International Cosmic Ray Conference (ICRC2023)-Cosmic-Ray Physics (Indirect, CRI), Nagoya, Japan, 26 July–3 August 2023; p. 291.
33. Schlueter, F.A. Expected Sensitivity of the AugerPrime Radio Detector to the Masses of Ultra-High-Energy Cosmic Rays Using Inclined Air Showers. Ph.D. Thesis, Karlsruher Institut für Technologie (KIT), Karlsruhe, Germany, 2022.
34. Corstanje, A.; Buitink, S.; Desmet, M.; Falcke, H.; Hare, B.M.; Hörandel, J.R.; Huege, T.; Jhansi, V.B.; Karastathis, N.; Krampah, G.K.; et al. A high-precision interpolation method for pulsed radio signals from cosmic-ray air showers. *J. Instrum.* **2023**, *18*, P09005. [[CrossRef](#)]
35. scikit-learn. Available online: <https://scikit-learn.org/stable/> (accessed on 12 July 2024).
36. Chicco, D.; Jurman, G. The advantages of the Matthews correlation coefficient (MCC) over F1 score and Accuracy in binary classification evaluation. *BMC Genom.* **2020**, *21*, 6. [[CrossRef](#)]
37. The Pierre Auger collaboration. Deep-learning based reconstruction of the shower maximum  $X_{\max}$  using the water-Cherenkov detectors of the Pierre Auger Observatory. *J. Instrum.* **2021**, *16*, P07019. [[CrossRef](#)]

- 
38. ResNet-18. Available online: <https://pytorch.org/vision/master/models/generated/torchvision.models.resnet18.html> (accessed on 14 June 2024).
  39. Ramzan, F.; Khan, M.U.G.; Rehmat, A.; Iqbal, S.; Saba, T.; Rehman, A.; Mehmood, Z. A deep learning approach for automated diagnosis and multi-class classification of Alzheimer's disease stages using resting-state fMRI and residual neural networks. *J. Med. Syst.* **2020**, *44*, 37. [[CrossRef](#)]

**Disclaimer/Publisher's Note:** The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.