



UNIVERSITY OF LEEDS

This is a repository copy of *Explainable deep learning for automatic rock classification*.

White Rose Research Online URL for this paper:

<https://eprints.whiterose.ac.uk/207266/>

Version: Accepted Version

Article:

Zheng, D., Zhong, H. orcid.org/0000-0003-3889-4065, Camps-Valls, G. et al. (6 more authors) (2024) Explainable deep learning for automatic rock classification. *Computers & Geosciences*, 184. 105511. ISSN 0098-3004

<https://doi.org/10.1016/j.cageo.2023.105511>

© 2023, Elsevier. This manuscript version is made available under the CC-BY-NC-ND 4.0 license <http://creativecommons.org/licenses/by-nc-nd/4.0/>.

Reuse

This article is distributed under the terms of the Creative Commons Attribution-NonCommercial-NoDerivs (CC BY-NC-ND) licence. This licence only allows you to download this work and share it with others as long as you credit the authors, but you can't change the article in any way or use it commercially. More information and the full terms of the licence here: <https://creativecommons.org/licenses/>

Takedown

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing eprints@whiterose.ac.uk including the URL of the record and the reason for the withdrawal request.



eprints@whiterose.ac.uk
<https://eprints.whiterose.ac.uk/>

1 **Explainable Deep Learning for Automatic Rock Classification**

2
3 Dongyu Zheng^a, Hanting Zhong^{a*}, Gustau Camps-Valls^b, Zhisong Cao^c, Xiaogang Ma^d,
4 Benjamin Mills^e, Xiumian Hu^f, Mingcai Hou^a, Chao Ma^a

5
6 ^a *State Key Laboratory of Oil and Gas Reservoir Geology and Exploitation and Key*
7 *Laboratory of Deep-time Geography and Environment Reconstruction and Applications,*
8 *MNR, Chengdu University of Technology, Chengdu, 610059, China*

9 ^b *Image Processing Laboratory, Universitat de València, Paterna, 46980, Spain*

10 ^c *College of Computer Science and Cyber Security, Chengdu University of Technology,*
11 *Chengdu, 610059, China*

12 ^d *Department of Computer Science, University of Idaho, Moscow, ID 83844, USA*

13 ^e *School of Earth and Environment, University of Leeds, Leeds, LS2 9JT, UK*

14 ^f *School of Earth Sciences and Engineering, Nanjing University, Nanjing, 210023, China*

15 * *Correspondence: zhonghanting@cdut.edu.cn*

16
17
18
19
20
21
22
23

24 **ARTICLE INFO**

25 **Keywords:**

26 Explainable deep learning

27 Knowledge-infused machine learning

28 Model interpretability

29 Attention-based modal network

30 Rock classification

31

32 **Authorship contribution statement**

33 **DZ**, Conceptualization, Methodology, Writing & Editing - Original Draft; **HZ**, Formal

34 analysis, Investigation, Writing - Review & Editing; **GC**, Writing - Review & Editing; **ZC**,

35 Methodology, Software, Formal analysis, Writing - Original Draft; **XM**, Writing - Review &

36 Editing; **BM**, Writing - Review & Editing; **XH**, Writing - Review & Editing; **MH**, Supervision,

37 Funding acquisition; **CM**, Project administration, Funding acquisition

38

39

40

41

42

43

44

45

46

47

Abstract

48 As deep learning (DL) gains popularity for its ability to make accurate predictions in
49 various fields, its applications in geosciences are also on the rise. Many studies focus on
50 achieving high accuracy in DL models by selecting models, developing more complex
51 architectures, and tuning hyperparameters. However, the interpretability of these models,
52 or the ability to understand how they make their predictions, is less frequently discussed.
53 To address the challenge of high accuracy but low interpretability of DL models in
54 geosciences, we study rock classification from thin-section photomicrographs of six types
55 of sedimentary rocks, including quartz arenite, feldspathic arenite, lithic arenite, siltstone,
56 oolitic packstone, and dolomite. These rocks' characteristic framework grains and grain
57 textures are their distinguishing features, such as the rounded or oval ooids in oolitic
58 packstone. We first train regular DL models, such as ResNet-50, on these
59 photomicrographs and achieve an accuracy of over 0.94. However, these models make
60 classifications based on features like cracks, cements, and scale bars, which are irrelevant
61 for distinguishing sedimentary rocks in real-world applications. We then propose an
62 attention-based dual network incorporating both global (overall photomicrograph) and local
63 (distinguishing framework grains) features to address this issue. Our proposed model
64 achieves not only high accuracy (0.99) but also provides interpretable feature extractions.
65 Our study highlights the need to consider interpretability and geological knowledge in
66 developing DL models, in addition to aiming for high accuracy.

67 **Keywords:** Explainable deep learning; Knowledge-infused machine learning; Model
68 interpretability; Attention-based modal network; Rock classification

69 **1. Introduction**

70 Deep learning (DL) has been highly effective in a range of tasks in geosciences,
71 including capturing complex relationships in datasets, creating automatic analysis
72 pipelines and building large and efficient models for numerical simulation and inversion
73 (Bergen et al., 2019; Reichstein et al., 2019; Camps-Valls et al., 2021). This is partly due
74 to the large number of tunable parameters and nonlinear model structures available in DL
75 approaches. However, this over-parameterization also leads to reduced interpretability,
76 making it difficult for users to understand the results obtained with these methods
77 (Castelvecchi, 2016; Buhrmester et al., 2021). The lack of interpretability limits the
78 reliability and application of DL models (Mamalakis et al., 2022). Scientists can neither
79 verify whether the predictions of DL models are made based on reasonable references nor
80 can they improve the models' ability of generalization (e.g., Ebert-Uphoff and Hilburn,
81 2020).

82 Since the successful application of deep convolutional neural networks (DCNNs) to
83 the classification of photographs from a dataset of 1.2 million images with one thousand
84 classes (Krizhevsky et al., 2012), there have been numerous attempts to use DCNNs for
85 fossil classification (Romero et al., 2020; Liu et al., 2020, 2022) and mineral classification
86 (Maitre et al., 2019; Hao et al., 2019; Wang et al., 2021; Ge et al., 2021; Zheng et al., 2022)
87 from photomicrographs, cathodoluminescence and scanning electron microscope images
88 in the geoscience community. These efforts have largely focused on evaluating the
89 performance of DCNNs using metrics such as accuracy or mean average precision.
90 However, these numerical values can be sensitive to small changes in the input images,

91 and models with "high accuracy" may not necessarily be robust if they rely on irrelevant
92 features for classification (Lei et al., 2018; Yang et al., 2020). Given the inherent
93 heterogeneity of rocks and other geological objects, it is important to understand how
94 DCNNs make their classifications to ensure their applicability in real-world scenarios. To
95 address this issue, recent advances in interpretability algorithms, such as the Class
96 Activation Mapping (CAM) technique, can provide valuable insights into the decision-
97 making process of DCNNs by highlighting the specific regions of an image that are
98 responsible for the classification using gradient information (Zhou et al., 2016; Selvaraju et
99 al., 2017). While these algorithms have been demonstrated to be useful, they have not yet
100 been widely applied to geoscientific tasks. To develop robust and reliable DL models for
101 geosciences, it is necessary to incorporate interpretability algorithms to deduce the
102 decision-making process of DL models.

103 Rock classification is a fundamental task in geoscience that involves identifying rock
104 types based on observing framework grains, minerals, texture, and structures. The
105 traditional approach for studying features of rocks in detail consists of first slicing and then
106 mounting, which makes rock samples sliced into roughly 30-micrometers-thick thin
107 sections and mounted on glass slides. Then, for affordability and efficiency, thin sections
108 are usually prepared and photographed as three-channel digital images (red, green, and
109 blue, RGB), also known as photomicrographs. Geoscientists can observe thin sections
110 under polarized light microscopes or examine photomicrographs to observe rock
111 compositions, texture, and other characteristics. Sedimentary rocks cover approximately
112 three-quarters of the Earth's surface, and understanding sedimentary rock types is

113 important for characterizing the Earth's landscape and life over time, as well as for
114 assessing reservoir quality in the oil and gas industry (Dickinson and Suczek, 1979;
115 Garzanti et al., 2007; Boggs, 2009). As a result, photomicrograph examination has become
116 a standard workflow in sedimentary geology, and there have been many attempts to use
117 deep learning approaches to classify rocks based on photomicrographs (de Lima et al.,
118 2019; Koeshidayatullah et al., 2020; Tang et al., 2020; Su et al., 2020; Saxena et al., 2021;
119 Li et al., 2022; Liu et al., 2022). These studies have primarily focused on the accuracy of
120 the models but have not yet investigated how deep learning models make their
121 classifications, which may lead to issues with generalization. This high accuracy yet low
122 interpretability of DL models for geosciences restricts the utility of DL in real-world
123 geoscience work.

124 In this study, we develop an interpretable rock classification DL model to address this
125 issue by incorporating geological knowledge. We focus on sedimentary rock classification
126 from thin-section photomicrographs, a common classification task in computer vision, as
127 the distinguishing features of sedimentary rocks, such as framework grains and textures,
128 are easy to identify visually. We first applied classical DCNNs, such as ResNet-50, to
129 classify six types of sedimentary rocks and evaluated the performance of these models
130 using numerical metrics (accuracy) and interpretable visualizations generated by Gradient-
131 weighted Class Activation Mapping (Grad-CAM). We then develop and test our new
132 attention-based dual-modal network, SedNet, which integrates global (the whole
133 photomicrograph) and local (characteristic framework grains) features. Our results show
134 that classical DCNNs achieve high accuracy but tend to focus on irrelevant parts of the

135 rock photomicrographs, while our proposed model achieves not only high accuracy but
136 also better interpretability, as indicated by the highlighting of distinguishing framework
137 grains in the Grad-CAM visualizations. This study underscores the importance of
138 interpretability and incorporation of geological knowledge in DL geoscience models. It
139 suggests that integrating global and local information may improve the generalization
140 abilities of DL models in this field.

141

142 **2. Dataset, data preprocessing, and data augmentation**

143 Six types of sedimentary rocks, including quartz arenite, feldspathic arenite, lithic
144 arenite, siltstone, oolitic packstone, and dolomite, were selected in this study for
145 classification (Figure 1; Table 1). Quartz arenite, feldspathic arenite, lithic arenite, and
146 siltstone are four types of siliciclastic rocks, consisting of grains formed by the
147 decomposition of rocks following weathering and deposition. These grains are typically
148 silicates such as quartz and feldspar, but may also include fragments of igneous,
149 metamorphic, and sedimentary rocks.

150 The differentiation between siliciclastic rocks lies in the composition of the framework
151 grains and grain size. For instance, quartz arenite is a sandstone with more than 95%
152 quartz grains. Therefore, a representative sub-image for quartz arenite showcases a
153 typical quartz grain. Feldspathic arenite and lithic arenite resemble quartz arenite but
154 contain predominantly feldspar and lithic grains, respectively. As such, their corresponding
155 sub-images feature a feldspar grain and a lithic fragment, respectively. Siltstone is a rock
156 type with smaller grains than sandstone, typically less than 0.063 mm, presenting a distinct

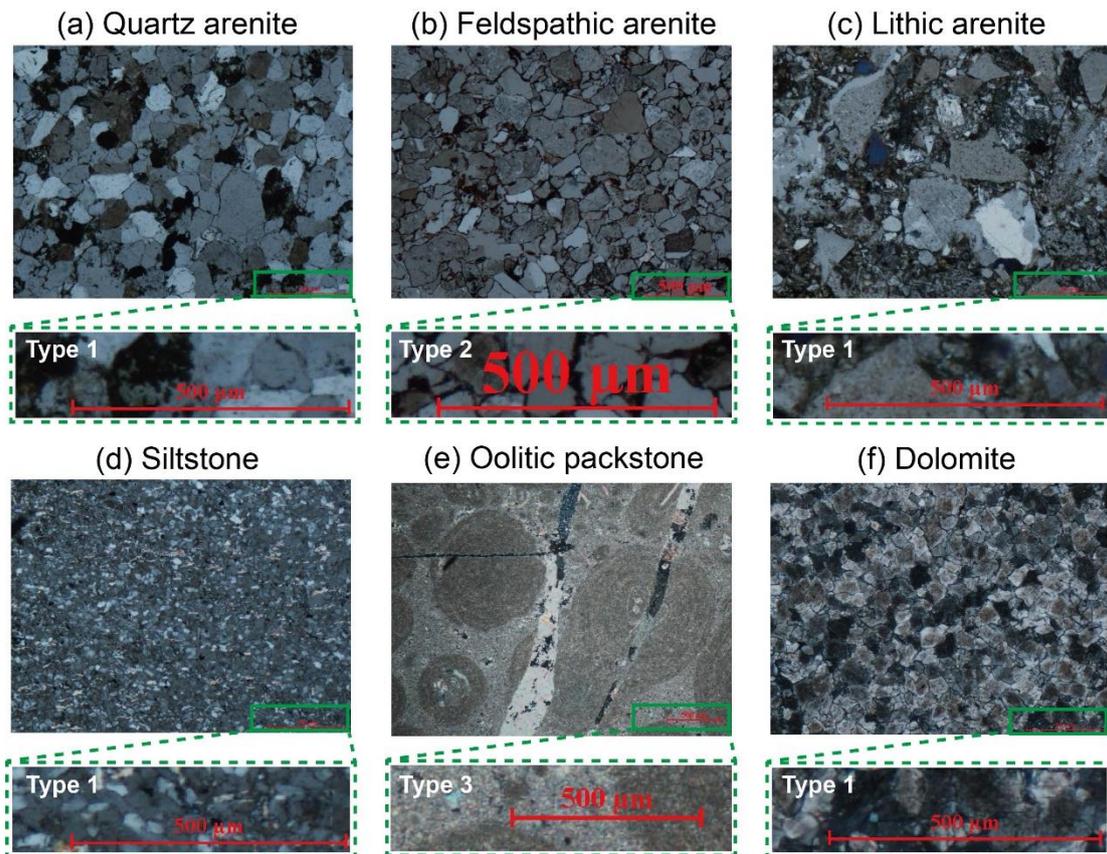
157 silty texture. A region displaying this texture was chosen as the siltstone sub-image, rather
158 than a single mineral grain.

159 In contrast to siliciclastic rocks, carbonate rocks like oolitic packstone and dolomite
160 are formed through chemical or biochemical processes, consisting primarily of calcite or
161 dolostone. Oolitic packstone, characterized by ooids—rounded or oval grains with
162 concentric textures—is represented by a sub-image showcasing an ooid. Dolomite, known
163 for its high interference colors and euhedral crystal forms, is represented by a sub-image
164 highlighting these unique attributes. In this manner, the sub-images for each rock type have
165 been carefully selected to encapsulate their unique mineralogical and textural
166 characteristics, thereby aiding the deep learning network in differentiating between the
167 classes more effectively.

168

169

170



171
172 Figure 1. The studied six types of sedimentary rocks and the associated scale bars.

173

174 A total of 1356 cross-polarized light photomicrographs were obtained from 15 samples

175 covering the six rock types using high-resolution electronic cameras mounted on Nikon

176 LV100POL microscopes. The images are three-channel RGB with a resolution of $1280 \times$

177 860×3 . The image dataset was split into training, validation, and test sets with a ratio of

178 6:2:2 (Figure 2). The images were subjected to flipping and rotation to augment the data

179 for model training. The color information in the images was considered important for

180 mineral differentiation and was therefore preserved. Furthermore, the overrepresentation

181 of certain rock types in the dataset (namely quartz arenite, feldspathic arenite, lithic arenite,

182 and dolomite) as opposed to others (siltstone and oolitic packstone) was motivated by the

183 inherent complexity associated with distinguishing these rock types. Quartz arenite,

184 feldspathic arenite, lithic arenite, and dolomite tend to exhibit homogenous grain
 185 compositions and textures, thereby posing a significant challenge in their identification.
 186 Conversely, rock types such as siltstone and oolitic packstone demonstrate distinct
 187 characteristics, such as fine-sized silty textures and oolitic grains, respectively. This
 188 deliberate imbalance in the dataset was designed to accommodate these differential
 189 complexities inherent in rock identification.

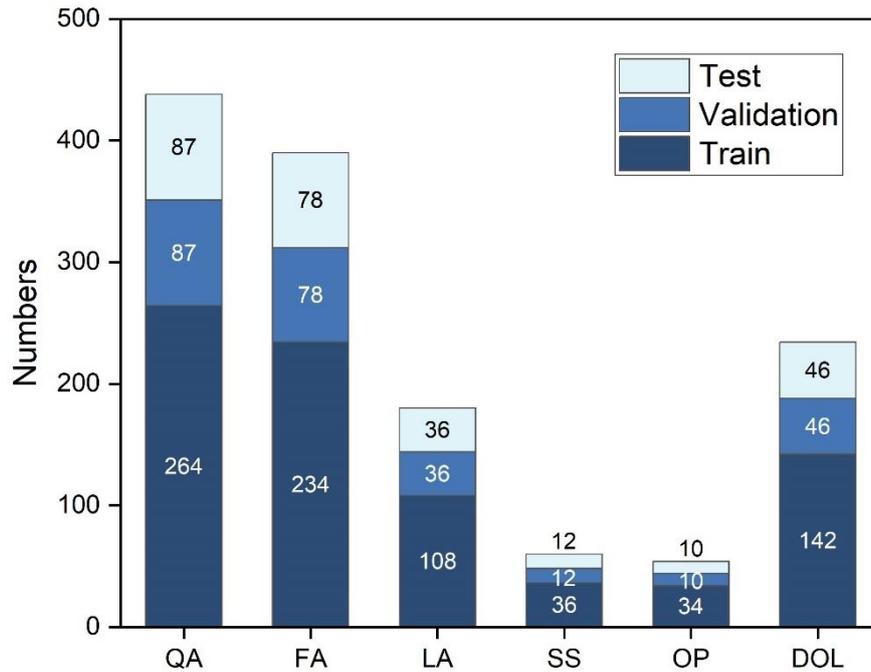
190

Table 1. Rock type descriptions

Rock type	Framework grains	Grain size (mm)	Distinguishing features	Scalar bar type
Quartz arenite	>90% quartz grains, trace feldspar or lithic fragments	~ 0.25-0.5	Over 90% quartz grains that are gray and clean under XPL	Type 1
Feldspathic arenite	~40-60% quartz grains, 20-40% feldspar grains, the rest are lithic grains	~ 0.25-0.5	Over 1/5 feldspar grains that are grey and mostly dirty. Tabular in shape, with cleavages or twinning under XPL	Type 2
Lithic arenite	30-50% lithic grains, the rest are quartz and feldspar grains	~ 0.063-0.5	Over 1/3 lithic grains that are volcanic or sedimentary rock fragments	Type 1
Siltstone	Mostly are quartz and feldspar grains	<0.063	Silty-sized grains that are grey under XPL	Type 1
Oolitic packstone	Ooids	0.25-2	Oval or rounded ooids grains with concentric fabrics	Type 2
Dolomite	Dolostone	<0.063	Euhedral-subhedral dolostone that are colorful under XPL	Type 1

See scalar bar type in Figure 1.

191



192

193 Figure 2. Numbers of training, validation, and test datasets. QA, quartz arenite; FA,
 194 feldspathic arenite; LA, lithic arenite; SS, siltstone; OP, oolitic packstone; DOL, dolomite.

195

196 3. Model implementation and interpretability

197 3.1 Model architecture

198 We introduce a dual-modal network called SedNet to classify sedimentary rocks using
 199 thin-section photomicrographs. Traditional DCNNs can only focus on a small area in the
 200 images due to the limited size of the convolution kernels. Therefore, SedNet incorporates
 201 an attention mechanism and dual-modal input to obtain both kernel-sized and grain-sized
 202 information, resulting in improved performance. The architecture of SedNet is depicted in
 203 Figure 3a. It consists of four modules: a dual feature extraction module, a channel attention
 204 module, a fused feature extraction module, and an output module. The dual feature
 205 extraction module comprises two parallel CNNs, each with two Residual Convolution
 206 blocks and a global pooling layer (Conv blocks in Figure 3a) activated by Rectified Linear
 207 Unit (ReLU) which returns the input values if the input is positive, and 0 if the input is

208 negative. To capture both global and local features in the rock classification, the dual
 209 feature extraction module takes as input both the rock photomicrographs and cropped
 210 images of distinctive framework grains within the original photomicrographs. For instance,
 211 for feldspathic arenite images, one input would be the original thin-section images, while
 212 the other would be a representative euhedral-subhedral feldspar grain with twining (see
 213 Section 2 for details). The channel attention module includes Squeeze-and-Excitation (SE)
 214 blocks, adapted from SE-Net (Hu et al., 2018). These SE blocks enable the neural network
 215 to emphasize important features and suppress less important features of the input data.
 216 This assignment procedure of feature importance is achieved through the Squeeze and
 217 Excitation operations. The Squeeze operation is a global pooling operation that converts a
 218 $N \times N \times C$ matrix (U_c) into a $1 \times 1 \times C$ matrix (Z_c), as shown in Eq. (1):

$$219 \quad Z_c = F_{sq}(U_c(i, j)) = \frac{1}{N^2} \sum_{i=1}^N \sum_{j=1}^N U_c(i, j), (1)$$

220 Where F_{sq} represents the Squeeze operation, c represents c^{th} channel, i and j
 221 represent the element of the i^{th} row and j^{th} column, respectively. The output of the
 222 Squeeze operation is a vector that contains information of the input feature maps (Squeeze
 223 vector in Figure 3a).

224 The Excitation operation takes the output of the Squeeze operation as the input, and
 225 then produces a set of channel-wise weights for each feature. This operation processes
 226 the $1 \times 1 \times C$ matrix with two fully connected layers and applies the sigmoid activation
 227 function to limit the output values to the range between 0 and 1. These values are then
 228 multiplied by the original $N \times N \times C$ matrix for each channel, as shown in the following
 229 equations:

230
$$\alpha = F_{\text{ex}}(Z, W) = \sigma(g(Z, W)) = \sigma(W_2 \delta(W_1 Z)), (2)$$

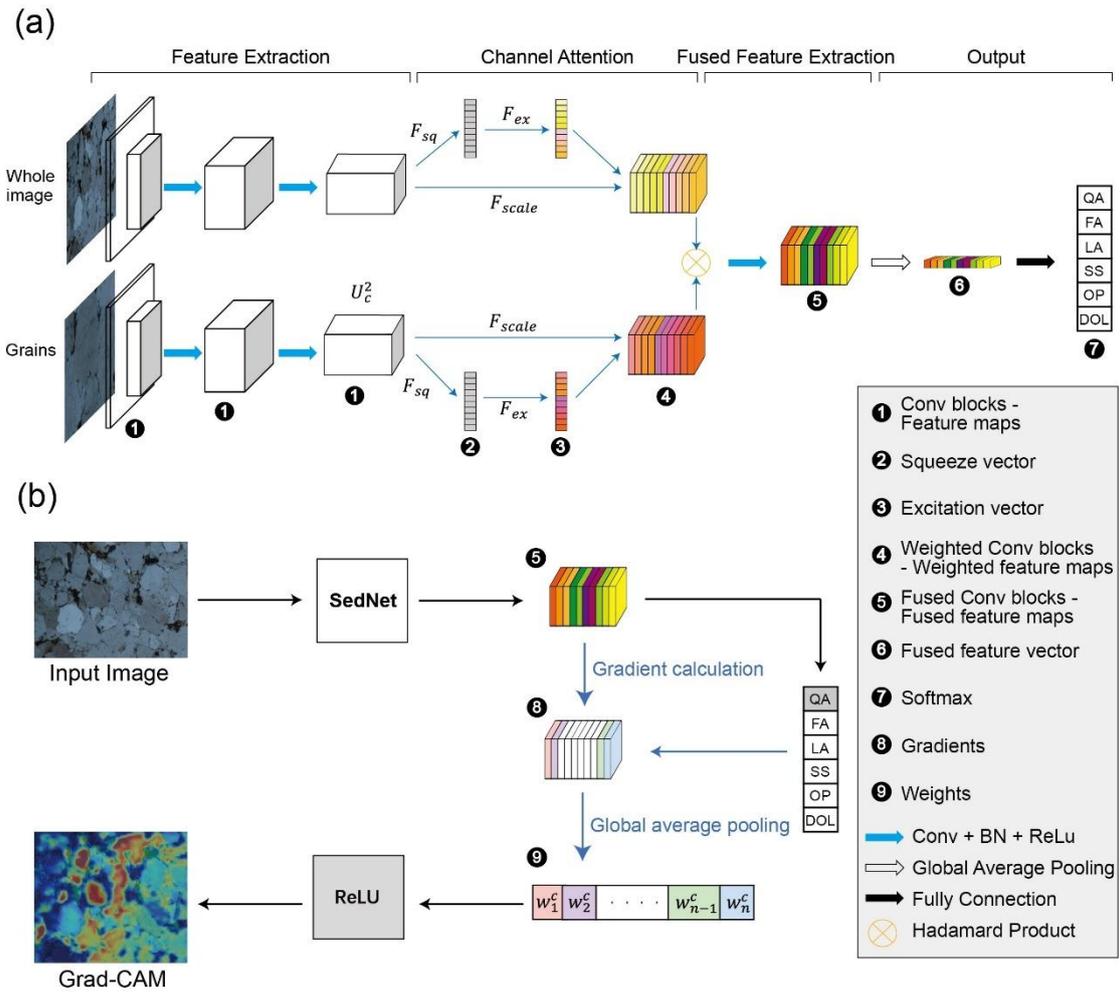
231
$$X_c = F_{\text{scale}}(U_c, \alpha_c) = \alpha_c U_c, (3)$$

232 α is a vector representing the weight of each channel, W_1 and W_2 are learnable weight
 233 matrices. g represents the dimensionality-reduction procedure, δ denotes the ReLU
 234 activation function, σ denotes the sigmoid activation function, X_c is the output of the
 235 channel attention module, α_c is the c^{th} element of α . The output of the Excitation
 236 operation contains information from the input feature maps with weighted feature
 237 importance (Excitation vector in Figure. 3a).

238 The fused feature extraction module operates the Hadamard product, also known as
 239 the element-wise product, and this operation enables the neural network to contain both
 240 global and local information from the input rock images. The Hadamard product takes two
 241 matrices of the same size and produces a new matrix where each element is the product
 242 of the corresponding elements of the original matrices (Fused Conv blocks in Figure 3a;
 243 as shown in Eq. 4). For example, given matrices A and B of dimensions $m \times n$ from the
 244 previous channel attention blocks, the fused feature extraction is the Hadamard product of
 245 A and B :

246
$$(A \odot B)_{ij} = (A)_{ij} * (B)_{ij}, (4)$$

247 \odot and $*$ denote the Hadamard operation, which is an element-wise multiplication, i and
 248 j represent the element of the i^{th} row and j^{th} column, respectively. The output module
 249 consists of global pooling (Fused feature vector in Figure 3a), fully connected and softmax
 250 layers (softmax in Figure 3a) that calculates the predicted probability for each category and
 251 outputs the category with the highest probability as the final prediction.



253

254

255

256

257

258

259 **3.2 Grad-CAM visualization**

260

261

262

263

264

Figure 3. (a) Overall structure of the SedNet. (b) Workflow of the calculation of Grad-CAM, adapted from Selvaraju et al. (2017). The marked numbers indicate key components of SedNet; see text for explanations.

Grad-CAM is a technique that visualizes the regions in an image that are most influential in the decision-making process of DCNNs (Selvaraju et al., 2017). In contrast to the traditional CAM algorithm, Grad-CAM offers enhanced performance due to its distinct approach to calculate class activation maps. Rather than employing the global average pooling over the feature map of the last convolution layer, as done in the CAM, Grad-CAM

265 computes the gradient of the output class with respect to the final convolutional layer of
266 the DCNN. This gradient information is then globally average-pooled to yield the neuron
267 importance weights. These weights are crucial in highlighting which features in the map
268 are most important for predicting the class, thereby resulting in a more effective and
269 comprehensive visualization of the class activations. In this study, we calculated the
270 gradient of the output of the DCNNs with respect to the feature maps of the last convolution
271 layer (Gradient in Figure 3b). The target convolutional layer is M^k , and k denotes the k^{th}
272 convolutional layer (Figure 3b). The weight of M^k can be calculated by Eq. (5):

$$273 \quad w_k^c = \frac{1}{N^2} \sum_{i=1}^N \sum_{j=1}^N \frac{\partial y^c}{\partial M_{ij}^k}, \quad (5)$$

274 Where c represents a category, w_k^c is a vector representing the weight of each channel
275 (Weights in Figure 3b) in M^k , y^c represents the score belonging to a certain category c .
276 Then the Grad-CAM can be obtained through the weighted combination of forward
277 activation and follow it with ReLU:

$$278 \quad I_k^c = \text{ReLU}(\sum_{k=1}^C w_k^c \cdot M^k), \quad (5)$$

279 Where I_k^c represents the Grad-CAM of the target convolution layer.

280

281 **4. Results and discussion**

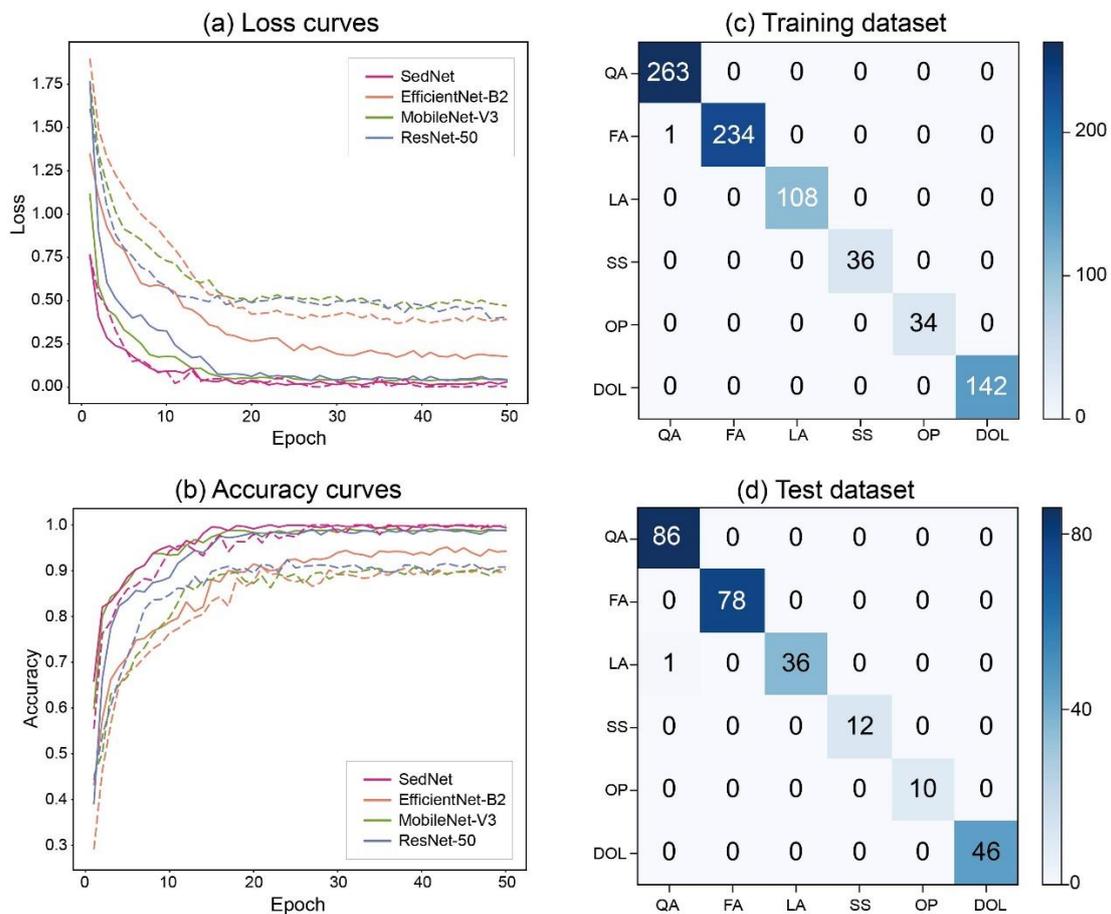
282 To thoroughly evaluate both the classical DCNNs and SedNet, we compared
283 numerical metrics and Grad-CAM visualizations. To ensure a fair comparison, we used the
284 same datasets and hyperparameters, such as training and test data, data augmentation,
285 batch size, and learning rate, for all models. Additionally, to test the interpretability of
286 DCNNs, we used scale bars for different rock types with various lengths and sizes. This

287 allowed us to better understand how the models were making their classifications.

288

289 4.1 The classical models have high accuracies but tend to use irrelevant feature 290 extractions

291 The three classical classification models, EfficientNet-B2 (Tan and Le, 2019),
292 MobileNet-V3 (Howard et al., 2017), and ResNet-50 (He et al., 2016), achieved high
293 accuracy and low loss values. ResNet-50 and MobileNet-V3 rapidly converged at
294 approximately 16 epochs, while EfficientNet-B2 reached its plateau after approximately 30
295 epochs of training (Figure 4). The training accuracies for EfficientNet-B2, MobileNet-V3,
296 and ResNet-50 are 0.9425, 0.9886, and 0.9871, while the test accuracies are 0.9187,
297 0.9458, and 0.9474 (Figure 4; Table 2).



298

299 Figure 4. (a) and (b) show the loss and accuracy of the EfficientNet-B2, MobileNet-V3,
 300 ResNet-50, and SedNet. The results of training dataset are represented in solid lines, and
 301 results of validation dataset are represented in dashed lines. (c) and (d) are confusion
 302 matrices of SedNet using the training and test dataset, respectively. QA, quartz arenite; FA,
 303 feldspathic arenite; LA, lithic arenite; SS, siltstone; OP, oolitic packstone; DOL, dolomite.

304

Table 2. Evaluation metrics and loss of trained deep learning models

		SedNet	EfficientNet- B2	MobileNet- V3	ResNet- 50
Accuracy	Train	0.9988	0.9360	0.9878	0.9866
	Test	0.9963	0.8993	0.9030	0.9104
	Mean	0.9975	0.9176	0.9454	0.9485
Precision	Train	1.0000	0.9651	0.9963	0.9926
	Test	1.0000	0.9414	0.9453	0.9457
	Mean	1.0000	0.9533	0.9708	0.9692
Recall	Train	0.9988	0.9688	0.9914	0.9938
	Test	0.9963	0.9526	0.9528	0.9606
	Mean	0.9975	0.9607	0.9721	0.9772
F1-score	Train	0.9994	0.9669	0.9938	0.9932
	Test	0.9981	0.9470	0.9490	0.9531
	Mean	0.9988	0.9569	0.9714	0.9732
Loss	Train	0.0315	0.1773	0.0453	0.0425
	Validation	0.0011	0.3918	0.4704	0.4058
	Mean	0.0163	0.2846	0.2579	0.2242

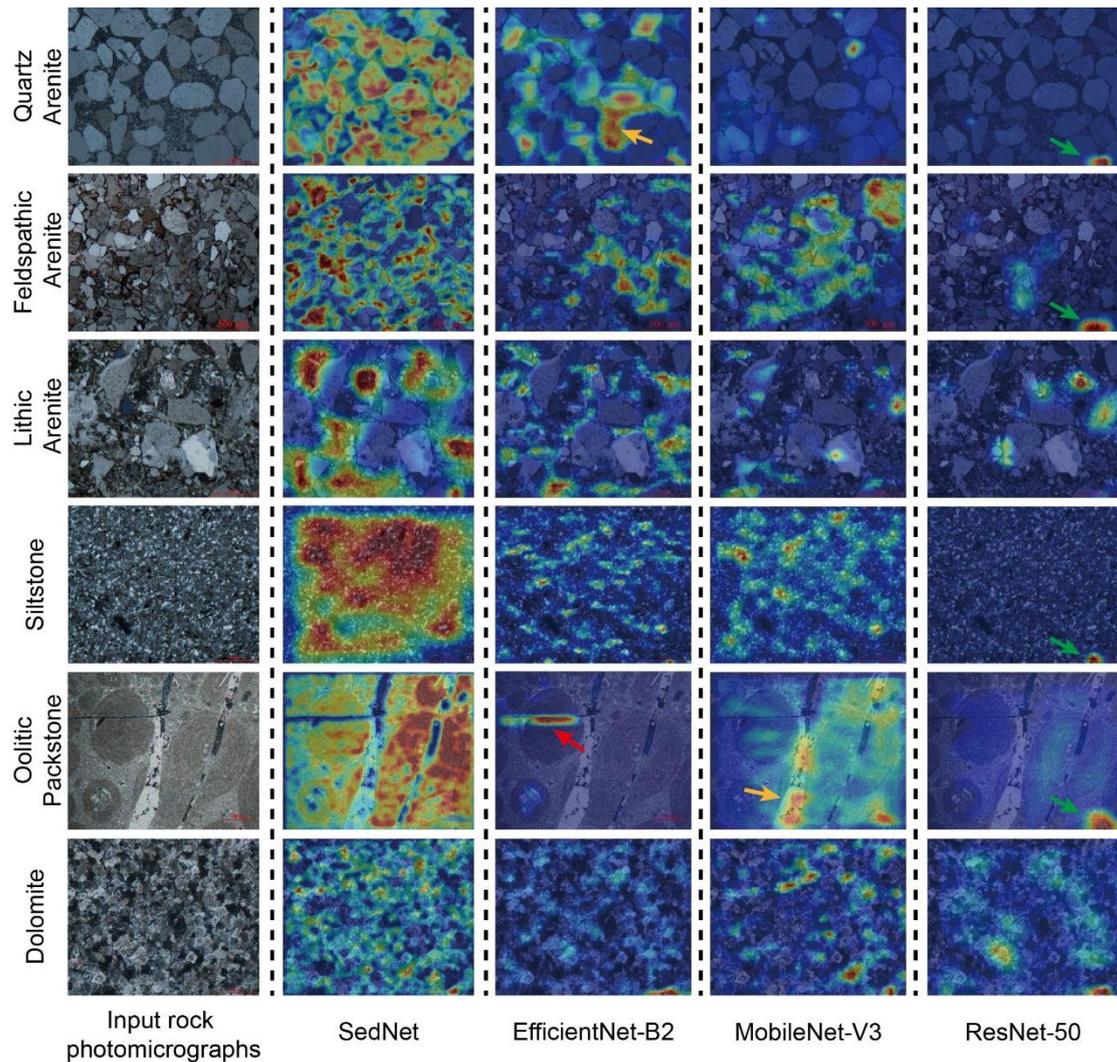
305

306 The Grad-CAM visualizations of the three models show uninterpretable features,
 307 suggesting that these models base their predictions on irrelevant features rather than the
 308 distinctive characteristics of sedimentary rocks (Figure 5). Although the characteristic
 309 features of the six rock types are distinct and easily recognizable by geologists, the
 310 classical DCNNs failed to focus on these features. For example, the characteristic features
 311 of quartz arenite, feldspathic arenite, lithic arenite and siltstone include the framework
 312 grains and their sandy or silty size, but the classical three models focus on other features
 313 such as the matrix and even scale bars. While the MobileNet-V3 model focuses on parts

314 of the framework grains in feldspar arenite, the visualization map indicates that it focuses
315 on the entire area rather than the distinctive feldspar grains.

316 Furthermore, the visualization map shows that the EfficientNet-B2 model identifies
317 siltstone based on a few grains. But, since the investigated siltstones are heterogeneous
318 in composition and grain size, the whole image or at least most of the siltstone
319 photomicrograph should be considered. The classical models also perform poorly with
320 carbonate rocks, where ooids, colorful interference colors, and euhedral dolostone crystal
321 forms are important features for geologists. However, the DCNNs again tend to focus on
322 the matrix, scale bar, and only a small number of grains (Figure 5). As the evaluation rules
323 of the classical DCNNs are not interpretable, these models may not be able to make
324 accurate predictions in real-world classification applications.

325



326

327 Figure 5. Grad-CAM of SedNet, EfficientNet-B2, MobileNet-V3, and ResNet-50. The red
 328 highlighted regions are the parts where models give more weight. Yellow arrows indicate
 329 the highlighted cements, red arrows indicate highlighted cracks, the green arrows indicate
 330 highlighted scale bars.

331

332 4.2 SedNet has not only high accuracy but also interpretable feature extractions

333 In comparison to the classical DCNNs, SedNet achieved not only high accuracy but
 334 also interpretable visualization maps. SedNet quickly converged at around 16 epochs, with
 335 training and test accuracies of 0.9942 and 1 (Figure 4; Table 2). The confusion matrix for
 336 SedNet shows that the model made very few mistakes, with only one feldspathic arenite
 337 being misclassified as a quartz arenite (Figure 4). Given the close similarity between quartz

338 and feldspathic arenite with high quartz content, SedNet demonstrated excellent
339 classification performance.

340 In addition to its accuracy, the Grad-CAM visualizations of SedNet are interpretable
341 and align with geologists' knowledge (Figure 5). The visualization maps for quartz,
342 feldspathic, and lithic arenite show that SedNet focuses on the most prominent quartz,
343 feldspar, and lithic grains, indicating that the classification was based on these distinctive
344 framework grains. For siltstone, the majority of the area in the photomicrographs is
345 highlighted, consistent with the distinguishing feature of silty texture. Similarly, the
346 visualization maps for the two carbonate rocks provide interpretable results with the ooids
347 and dolostone in the photomicrographs being highlighted. SedNet used the characteristic
348 framework grains and overall rock textures, rather than the matrix, scale bars, and cracks,
349 to classify the rocks, which mimics the approach of geologists and therefore has a high
350 potential for real rock classification projects.

351

352 **4.3 Maximizing accuracy and decision-making power: the importance of** 353 **interpretability and geological knowledge**

354 In geosciences, accuracy is a crucial factor in developing and evaluating DL models.
355 However, we argue that geological knowledge plays a key role in providing context and
356 understanding for accurate predictions, and therefore its incorporation into DL models is
357 necessary. Without this knowledge, even highly accurate models may not be reliable or
358 meaningful. Most contemporary CNNs (e.g., ResNet, He et al., 2016) are characterized by
359 millions of parameters, necessitating extensive image datasets for optimal performance.

360 Although these models are effective in numerous tasks, their suitability for specific domains
361 necessitates additional validation. For example, in tasks like rock image identification,
362 acquiring a large dataset is often unfeasible. Traditional CNNs tend to overfit by
363 memorizing distinct features (e.g., cracks, as mentioned in this study) in smaller datasets,
364 rather than identifying the features geologists typically use for rock classification. However,
365 this overfitting is often undetected when evaluated solely through the numerical metrics.
366 Therefore, the integration of geological knowledge into the development and evaluation of
367 DL models for geoscience applications is essential to ensure accuracy, relevance, and
368 interpretability (Barnes et al., 2020; McGovern et al., 2019; Ebert-Uphoff and Hilburn, 2020).

369 There is a trend in the geoscience community towards using machine learning models
370 with a better understanding of their inner workings. For example, recent research (Zhao et
371 al., 2019; Doucet et al., 2022; Zou et al., 2022) introduced feature importance and Shapley
372 Additive Explanations (SHAP; Lundberg and Lee, 2017), an interpretable algorithm based
373 on cooperative game theory, into geochemistry studies. Using SHAP values, trace
374 elements in basalts can be used to classify tectonic settings and identify new geochemical
375 differences between basalts from convergent and divergent boundaries. Toms et al. (2020)
376 introduced layerwise relevance propagation for identifying meaningful patterns in ENSO
377 (El Niño-Southern Oscillation) phase identification and seasonal prediction. This algorithm
378 provides transparency to machine learning by propagating relevance from the output layer
379 back through the network to the input layer based on the relative contribution of each
380 neuron (Bach et al., 2015). These interpretable algorithms have shown great value in
381 meteorology, being used to make subseasonal forecasts (Mayer and Barnes, 2021) and to

382 reveal slowdowns in decadal climate warming (Labe and Barnes, 2022).

383 In our study, geologists can easily distinguish sedimentary rocks based on their unique
384 framework grains and textures. However, classical convolutional neural networks (DCNNs)
385 often placed more weight on scale bars, cements, and cracks, as evident in the Gradient-
386 weighted Class Activation Mapping (Grad-CAM) visualizations (Figure 5). The rock
387 photomicrographs used in the study were intentionally presented with various styles of
388 scale bars (Figure 1) and the studied quartz arenite was the only rock type with cracks. As
389 a result, the scale bars and cracks became the most distinguishing features of the DCNNs.
390 Such "noise" can be introduced easily if the rock samples are photographed by different
391 institutes or in different facilities, and can significantly impact the final outputs. It has been
392 noted that even slight image transformations can alter the predictions of DCNNs (Azulay
393 and Weiss, 2018). One potential solution to this generalization issue is a collection of big
394 datasets. However, it may not always be feasible for geologists to collect the same volume
395 of images as datasets such as ImageNet. In these cases, the assessment of DCNN outputs
396 cannot rely solely on numerical metrics like accuracy, as these only indicate the algorithm's
397 performance on a known dataset. Without a full understanding of how the algorithm works,
398 trained models may still face generalization issues. Our proposed dual-modal network is a
399 potential solution for image classification tasks. It emphasizes distinguishing features in
400 the DL models and achieves high accuracy and interpretable Grad-CAM visualizations by
401 integrating global and local features (Figure 5).

402 The proposed dual-modal network represents one way in enhancing the performance
403 and interpretability of DCNNs, achieved by integrating domain-specific modules. In addition

404 to the integration of bespoke modules into the network, alternative strategies also show
405 potential in this respect. A major challenge in applying deep learning to geosciences is the
406 lack of labeled training data, often due to subjective or labor-intensive labeling processes.
407 One solution is to generate synthetic datasets based on fundamental geological principles.
408 This not only helps in tuning model parameters but also ensures model robustness when
409 applied to real-world problems. This technique has been used successfully in various
410 applications, such as detecting permeability from rock images by generating samples of
411 porous media and assigning permeability labels using the Boltzmann method (Wu et al.,
412 2018), and seismic interpretation where large datasets can be generated through forward
413 modelling (Wu et al., 2023).

414 Another approach is to integrate prior knowledge into the loss function. This requires
415 DCNNs to conform to the training dataset while also adhering to the prior knowledge, such
416 as physical laws defined by partial differential equations, and this type of neural network is
417 also known as the physics-informed neural network. The powerful approximation and high
418 expressivity capacities of DCNNs enable them to infer solutions within the complex space
419 defined by the governing physical laws, thereby enhancing model performance and
420 understanding (Cuomo et al., 2022; Zhang et al., 2023). In the realm of rock imaging, the
421 prior knowledge can be the experience of geologists that some rocks or grains are more
422 identifiable due to their distinctive features. By applying different weights to the
423 regularization terms in the loss function, the model's ability to analyze difficult images can
424 be improved. Given the desire for a physical understanding of DL models in geosciences,
425 it is expected that interpretability and geological knowledge will be further incorporated in

426 future applications of DCNNs. Combining interpretability and geological knowledge with
427 Deep Learning can create more effective and transparent models for geoscience
428 applications.

429

430 **5. Conclusions**

431 To examine the significance of interpretability in DL models for geoscientific tasks, we
432 conducted automatic sedimentary rock classification using thin-section photomicrographs
433 and DL models. While classical deep convolutional neural networks (DCNNs) such as
434 ResNet, EfficientNet, and MobileNet achieved high accuracy (up to 0.96), their Grad-CAM
435 visualizations were often not geological-reasonable. Framework grains and textures are
436 key features for distinguishing the studied sedimentary rocks. However, classical DCNNs
437 tended to classify rocks based on irrelevant features, as indicated by the highlighted
438 regions of scale bars, cements, and cracks. To address this issue, we proposed an
439 attention-based dual network that inputs both the original thin-section photomicrographs
440 and framework grains. By combining information from the whole images and framework
441 grains, our proposed model achieves not only higher accuracy (0.99) but also produces
442 interpretable visualization heatmaps in which framework grains were given more weight in
443 the classification process. Our study highlights the importance of considering
444 interpretability and geological knowledge in developing DL models, in addition to aiming
445 for high accuracy.

446

447 **Acknowledgement**

448 The authors would like to thank Tingting Gong for providing thin sections and Jiaqi Wu
449 for helping in capturing photomicrographs. Thanks to three anonymous reviewers for
450 constructive reviews that substantially improved the manuscript. This work was financially
451 supported by National Natural Science Foundation of China (Grant Nos. 41888101,
452 42050104, 42050102, 42202125, 42172137). Additional support was provided by the IUGS
453 Deep-time Digital Earth (DDE) Big Science Program.

454

455 **Code availability statement**

456 The code is made open on Github repository at:

457 https://github.com/MudRocw1/SedNet_explainable-deep-learning-network

458

459 **References**

460

461 Azulay, A., Weiss, Y., 2018. Why do deep convolutional networks generalize so poorly to
462 small image transformations? arXiv Prepr. arXiv1805.12177.

463 Bach, S., Binder, A., Montavon, G., Klauschen, F., Müller, K.-R., Samek, W., 2015. On pixel-
464 wise explanations for non-linear classifier decisions by layer-wise relevance
465 propagation. PLoS One 10, e0130140.

466 Barnes, E.A., Toms, B., Hurrell, J.W., Ebert-Uphoff, I., Anderson, C., Anderson, D., 2020.
467 Indicator Patterns of Forced Change Learned by an Artificial Neural Network. J. Adv.
468 Model. Earth Syst. 12. <https://doi.org/10.1029/2020MS002195>

469 Bergen, Karianne J., Paul A. Johnson, V. Maarten, Beroza, G.C., 2019. Machine learning for
470 data-driven discovery in solid Earth geoscience. Science (80-.). 363.
471 <https://doi.org/10.1126/science.aau0323>

472 Boggs Jr, S., Boggs, S., 2009. Petrology of sedimentary rocks. Cambridge university press.

473 Buhrmester, V., Münch, D., Arens, M., 2021. Analysis of Explainers of Black Box Deep Neural
474 Networks for Computer Vision: A Survey. Mach. Learn. Knowl. Extr. 3, 966–989.
475 <https://doi.org/10.3390/make3040048>

476 Camps-Valls, G., Tuia, D., Zhu, X.X., Reichstein, M., 2021. Deep learning for the Earth
477 Sciences: A comprehensive approach to remote sensing, climate science and
478 geosciences. John Wiley & Sons.

479 Castelvechi, D., 2016. The black box of AI. Nature 538, 20–23.

480 Cuomo, S., Di Cola, V.S., Giampaolo, F., Rozza, G., Raissi, M., Piccialli, F., 2022. Scientific
481 machine learning through physics-informed neural networks: Where we are and what's

482 next. *J. Sci. Comput.* 92, 88.

483 de Lima, R.P., Bonar, A., Duarte Coronado, D., Marfurt, K., Nicholson, C., 2019. Deep
484 convolutional neural networks as a geological image classification tool. *Sediment. Rec.*
485 17, 4–9. <https://doi.org/10.2110/sedred.2019.2.4>

486 Dickinson, W.R., Suczek, C.A., 1979. Plate tectonics and sandstone compositions. *Am.*
487 *Assoc. Pet. Geol. Bull.* 63, 2164–2182.

488 Doucet, L.S., Tetley, M.G., Li, Z.-X., Liu, Y., Gamaleldien, H., 2022. Geochemical
489 fingerprinting of continental and oceanic basalts: A machine learning approach. *Earth-*
490 *Science Rev.* 104192.

491 Ebert-Uphoff, I., Hilburn, K., 2020. Evaluation, tuning, and interpretation of neural networks for
492 working with images in meteorological applications. *Bull. Am. Meteorol. Soc.* 101,
493 E2149–E2170. <https://doi.org/10.1175/BAMS-D-20-0097.1>

494 Garzanti, E., Doglioni, C., Vezzoli, G., Ando, S., 2007. Orogenic belts and orogenic sediment
495 provenance. *J. Geol.* 115, 315–334.

496 Ge, S., Wang, C., Jiang, Z., Hao, H., Gu, Q., 2021. Dual-input attention network for automatic
497 identification of detritus from river sands. *Comput. Geosci.* 151, 104735.
498 <https://doi.org/10.1016/j.cageo.2021.104735>

499 Hao, H., Guo, R., Gu, Q., Hu, X., 2019. Machine learning application to automatically classify
500 heavy minerals in river sand by using SEM/EDS data. *Miner. Eng.* 143, 105899.
501 <https://doi.org/10.1016/j.mineng.2019.105899>

502 He, Kaiming, Xiangyu Zhang, Shaoqing Ren, Sun, J., 2016. Deep residual learning for image
503 recognition., in: *Proceedings of the IEEE Conference on Computer Vision and Pattern*
504 *Recognition.* pp. 770–778.

505 Howard, A.G., Zhu, M., Chen, B., Kalenichenko, D., Wang, W., Weyand, T., Andreetto, M.,
506 Adam, H., 2017. Mobilenets: Efficient convolutional neural networks for mobile vision
507 applications. *arXiv Prepr. arXiv1704.04861*.

508 Hu, J., Shen, L., Sun, G., 2018. Squeeze-and-excitation networks, in: *Proceedings of the*
509 *IEEE Conference on Computer Vision and Pattern Recognition.* pp. 7132–7141.

510 Koeshidayatullah, A., Morsilli, M., Lehrmann, D.J., Al-Ramadan, K., Payne, J.L., 2020. Fully
511 automated carbonate petrography using deep convolutional neural networks. *Mar. Pet.*
512 *Geol.* 122, 104687. <https://doi.org/10.1016/j.marpetgeo.2020.104687>

513 Krizhevsky, A., Sutskever, I. and Hinton, G.E., 2012. Imagenet classification with deep
514 convolutional neural networks. *Adv. Neural Inf. Process. Syst.* 25, 1097–1105.

515 Labe, Z.M., Barnes, E.A., 2022. Predicting Slowdowns in Decadal Climate Warming Trends
516 With Explainable Neural Networks. *Geophys. Res. Lett.* 49.
517 <https://doi.org/10.1029/2022GL098173>

518 Lei, S., Zhang, H., Wang, K., Su, Z., 2018. How training data affect the accuracy and
519 robustness of neural networks for image classification.

520 Li, D., Zhao, J., Ma, J., 2022. Experimental Studies on Rock Thin-Section Image
521 Classification by Deep Learning-Based Approaches.

522 Liu, T., Li, C., Liu, Zongbao, Zhang, K., Liu, F., Li, D., Zhang, Y., Liu, Zhigang, Liu, L., Huang,
523 J., 2022. Research on Image Identification Method of Rock Thin Slices in Tight Oil
524 Reservoirs Based on Mask R-CNN. *Energies* 15. <https://doi.org/10.3390/en15165818>

525 Liu, X., Jiang, S., Wu, R., Shu, W., Hou, J., Sun, Y., Sun, J., Chu, D., Wu, Y., Song, H., 2022.

526 Automatic taxonomic identification based on the Fossil Image Dataset (> 415,000
527 images) and deep convolutional neural networks. *Paleobiology* 1–22.

528 Lundberg, S.M., Lee, S.-I., 2017. A unified approach to interpreting model predictions. *Adv.*
529 *Neural Inf. Process. Syst.* 30.

530 Maitre, J., Bouchard, K., Bédard, L.P., 2019. Mineral grains recognition using computer vision
531 and machine learning. *Comput. Geosci.* 130, 84–93.
532 <https://doi.org/10.1016/j.cageo.2019.05.009>

533 Mamalakis, A., Ebert-Uphoff, I., Barnes, E.A., 2022. Neural network attribution methods for
534 problems in geoscience: A novel synthetic benchmark dataset. *Environ. Data Sci.* 1, 1–
535 17. <https://doi.org/10.1017/eds.2022.7>

536 Mayer, K.J., Barnes, E.A., 2021. Subseasonal Forecasts of Opportunity Identified by an
537 Explainable Neural Network. *Geophys. Res. Lett.* 48, 1–9.
538 <https://doi.org/10.1029/2020GL092092>

539 McGovern, A., Lagerquist, R., Gagne, D.J., Jergensen, G.E., Elmore, K.L., Homeyer, C.R.,
540 Smith, T., 2019. Making the black box more transparent: Understanding the physical
541 implications of machine learning. *Bull. Am. Meteorol. Soc.* 100, 2175–2199.
542 <https://doi.org/10.1175/BAMS-D-18-0195.1>

543 Reichstein, M., Camps-Valls, G., Stevens, B., Jung, M., Denzler, J., Carvalhais, N., Prabhat,
544 2019. Deep learning and process understanding for data-driven Earth system science.
545 *Nature* 566, 195–204. <https://doi.org/10.1038/s41586-019-0912-1>

546 Romero, I.C., Kong, S., Fowlkes, C.C., Jaramillo, C., Urban, M.A., Oboh-Ikuenobe, F.,
547 D’Apolito, C., Punyasena, S.W., 2020. Improving the taxonomy of fossil pollen using
548 convolutional neural networks and superresolution microscopy. *Proc. Natl. Acad. Sci. U.*
549 *S. A.* 117, 28496–28505. <https://doi.org/10.1073/pnas.2007324117>

550 Saxena, N., Day-Stirrat, R.J., Hows, A., Hofmann, R., 2021. Application of deep learning for
551 semantic segmentation of sandstone thin sections. *Comput. Geosci.* 152, 104778.
552 <https://doi.org/10.1016/j.cageo.2021.104778>

553 Selvaraju, R.R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., Batra, D., 2017. Grad-cam:
554 Visual explanations from deep networks via gradient-based localization, in: *Proceedings*
555 *of the IEEE International Conference on Computer Vision*. pp. 618–626.

556 Su, C., Xu, S., Zhu, K., Zhang, X., 2020. Rock classification in petrographic thin section
557 images based on concatenated convolutional neural networks. *Earth Sci. Informatics* 13,
558 1477–1484.

559 Su, D., Zhang, H., Chen, H., Yi, J., Chen, P.-Y., Gao, Y., 2018. Is Robustness the Cost of
560 Accuracy?--A Comprehensive Study on the Robustness of 18 Deep Image
561 Classification Models, in: *Proceedings of the European Conference on Computer Vision*
562 *(ECCV)*. pp. 631–648.

563 Tan, M., Le, Q., 2019. Efficientnet: Rethinking model scaling for convolutional neural
564 networks, in: *International Conference on Machine Learning*. PMLR, pp. 6105–6114.

565 Tang, D.G., Milliken, K.L., Spikes, K.T., 2020. Machine learning for point counting and
566 segmentation of arenite in thin section. *Mar. Pet. Geol.* 120.
567 <https://doi.org/10.1016/j.marpetgeo.2020.104518>

568 Toms, B.A., Barnes, E.A., Ebert-Uphoff, I., 2020. Physically Interpretable Neural Networks for
569 the Geosciences: Applications to Earth System Variability. *J. Adv. Model. Earth Syst.* 12,

570 1–20. <https://doi.org/10.1029/2019MS002002>
571 Wang, C., Ge, S., Jiang, Z., Hao, H., Gu, Q., 2021. Computers and Geosciences
572 SiamFuseNet : A pseudo-siamese network for detritus detection from polarized
573 microscopic images of river sands. *Comput. Geosci.* 156, 104912.
574 <https://doi.org/10.1016/j.cageo.2021.104912>
575 Wu, J., Yin, X., Xiao, H., 2018. Seeing permeability from images: fast prediction with
576 convolutional neural networks. *Sci. Bull.* 63, 1215–1222.
577 Wu, X., Ma, J., Si, X., Bi, Z., Yang, J., Gao, H., Xie, D., Guo, Z., Zhang, J., 2023. Sensing
578 prior constraints in deep neural networks for solving exploration geophysical problems.
579 *Proc. Natl. Acad. Sci. U. S. A.* 120, 1–12. <https://doi.org/10.1073/pnas.2219573120>
580 Yang, Y.-Y., Rashtchian, C., Zhang, H., Salakhutdinov, R.R., Chaudhuri, K., 2020. A closer
581 look at accuracy vs. robustness. *Adv. Neural Inf. Process. Syst.* 33, 8588–8601.
582 Zhang, Y., Zhu, X., Gao, J., 2023. Seismic inversion based on acoustic wave equations using
583 physics-informed neural network. *IEEE Trans. Geosci. Remote Sens.* 61, 1–11.

584 Zhao, Y., Zhang, Y., Geng, M., Jiang, J., Zou, X., 2019. Involvement of slab-derived fluid in
585 the generation of Cenozoic basalts in Northeast China inferred from machine learning.
586 *Geophys. Res. Lett.* 46, 5234–5242.

587 Zheng, D., Wu, S., Ma, C., Xiang, L., Hou, L., Chen, A., Hou, M., 2022. Zircon classification
588 from cathodoluminescence images using deep learning. *Geosci. Front.* 101436.
589 <https://doi.org/10.1016/j.gsf.2022.101436>
590 Zou, S., Chen, X., Brzozowski, M.J., Leng, C., Xu, D., 2022. Application of machine learning
591 to characterizing magma fertility in porphyry Cu deposits. *J. Geophys. Res. Solid Earth*
592 127, e2022JB024584.

593
594
595

596 **List of Figures**

597
598
599

Figure 1. The studied six types of sedimentary rocks and the associated scale bars.

600
601
602

Figure 2. Numbers of training, validation, and test datasets. QA, quartz arenite; FA, feldspathic arenite; LA, lithic arenite; SS, siltstone; OP, oolitic packstone; DOL, dolomite.

603
604
605

Figure 3. (a) Overall structure of the SedNet. (b) Workflow of the calculation of Grad-CAM, adapted from Selvaraju et al. (2017). The marked numbers indicate key components of SedNet, see text for explanations.

606
607
608

Figure 4. (a) and (b) show the loss and accuracy of the EfficientNet-B2, MobileNet-V3, ResNet-50, and SedNet. (c) and (d) are confusion matrices of SedNet using the training and test dataset, respectively. QA, quartz arenite; FA, feldspathic arenite; LA, lithic arenite; SS, siltstone; OP, oolitic packstone; DOL, dolomite.

609
610
611

Figure 5. Grad-CAM of the SedNet, EfficientNet-B2, MobileNet-V3, and ResNet-50. The

612

613 red highlighted regions are the parts where models give more weight. Yellow arrows
614 indicate the highlighted cements, red arrows indicate highlighted cracks, the green arrows
615 indicate highlighted scale bars.

616

617 **List of Tables**

618

619 Table 1. Rock type descriptions

620

621 Table 2. Evaluation metrics and loss of the trained deep learning models