

1 **Revision 1**

2 **The application of “transfer learning” in optical microscopy: the petrographic**  
3 **classification of metallic minerals**

4 Word Count: 5753

5 **Yi-Wei Cai<sup>1</sup>, Kun-Feng Qiu<sup>1\*</sup>, Maurizio Petrelli<sup>2</sup>, Zhao-Liang Hou<sup>3</sup>, M. Santosh<sup>1,4</sup>, Hao-**  
6 **Cheng Yu<sup>1</sup>, Ryan T. Armstrong<sup>5</sup>, Jun Deng<sup>1,6</sup>**

7 *<sup>1</sup> Frontiers Science Center for Deep-time Digital Earth, State Key Laboratory of Geological*  
8 *Processes and Mineral Resources, School of Earth Sciences and Resources, China University of*  
9 *Geosciences, Beijing 100083, China*

10 *<sup>2</sup> Department of Physics and Geology, University of Perugia, Perugia 06100, Italy*

11 *<sup>3</sup> Department of Geology, University of Vienna, Vienna 1090, Austria*

12 *<sup>4</sup> Department of Earth Science, University of Adelaide, Adelaide SA 5005, Australia*

13 *<sup>5</sup> School of Minerals and Energy Resources Engineering, University of New South Wales, Sydney,*  
14 *NSW 2052, Australia*

15 *<sup>6</sup> Geological Research Institute of Shandong Gold Group Co., Ltd., Jinan 250013, China*

16

17 **\*Corresponding author**

18 Kun-Feng Qiu [kunfengqiu@qq.com](mailto:kunfengqiu@qq.com)

19 Professor, China University of Geosciences, Beijing

20 No. 29 Xueyuan Road, Haidian District, Beijing, 100083, P.R. China

21

## Abstract

22 Analysis of optical microscopic image data is crucial for the identification and  
23 characterization of mineral phases, and thus directly relevant to the subsequent methodology  
24 selections of further detailed petrological exploration. Here we present a novel application of Swin  
25 Transformer, a deep learning algorithm to classify metal mineral phases such as arsenopyrite,  
26 chalcopyrite, gold, pyrite, and stibnite, in images captured by optical microscopy. To speed up the  
27 training process and improve the generalization capabilities of the investigated model, we adopt  
28 the “transfer learning” paradigm by pretraining the algorithm using a large, general-purpose, image  
29 dataset named ImageNet-1k. Further, we compare the performances of the Swin Transformer with  
30 those of two well-established Convolutional Neural Networks (CNNs) named MobileNetv2 and  
31 ResNet50, respectively. Our results highlight a maximum accuracy of 0.92 for the Swin  
32 Transformer, outperforming the CNNs. To provide an interpretation of the trained models, we  
33 apply the so-called Class Activation Map (CAM), which points to a strong global feature extraction  
34 ability of the Swin Transformer metal mineral classifier that focuses on distinctive (e.g., colors)  
35 and microstructural (e.g., edge shapes) features. The results demonstrate that the deep learning  
36 approach can accurately extract all available attributes, which reveals the potential to assist in data  
37 exploration and provides an opportunity to carry out spatial quantization at a large scale (cm-mm).  
38 Simultaneously, boosting the learning processes with pre-trained weights can accurately capture  
39 relevant attributes in mineral classification, revealing the potential for application in mineralogy  
40 and petrology, as well as enabling its use in resource explorations.

41

42 **Keywords:** Swin Transformer metal mineral classifier; Microscopy images; Transfer learning;  
43 Deep learning; Class Activation Map

44

## Introduction

45 Petrographic studies at the microscopic scale and mineral identification constitute the  
46 fundamental step in many geological studies (e.g., igneous, metamorphic and sedimentary  
47 petrology or mineral exploration) and industrial productions (Schrader and Zega 2019; Deng et al.  
48 2020a; dos Anjos et al. 2021; Sheldrake and Higgins 2021; Azeuda Ndonfack et al. 2022). In  
49 petrographic investigations at microscopic scale, the first step mostly relies on optical microscopy  
50 involving the identification of mineral phases and textures (Su et al. 2020; Leichter et al. 2022;  
51 Faria et al. 2022; Qiu et al. 2023c). In recent years, optical microscopic observations are further  
52 supported by more advanced techniques, like electron-based imaging and X-ray techniques (Fu  
53 and Aldrich 2019). Meanwhile, many software like ImageJ (Schneider et al. 2012) or scripting  
54 languages, such as Python (Petrelli 2021) or MATLAB (Trauth et al. 2007) now support  
55 quantitative petrographic investigations like the segmentation processes or crystal size distribution  
56 analyses (Santosh et al. 2009; Tarquini and Favalli 2010; Jungmann et al. 2014; Y. Wang et al.  
57 2021; Zhang et al. 2021). Despite the recent analytical advancements and the possibilities of  
58 automation in quantitative petrographic studies, the initial investigation of new samples still relies,  
59 mostly, on the manual identification of mineral phases by expert petrologists, by optical  
60 microscopy (Deng et al. 2020b). This procedure is time-consuming, often subjective, and  
61 sometimes biased since many minerals share similar textural and optical properties (Santosh 2010;  
62 Młynarczyk et al. 2013; Xu et al. 2021; Zhong et al. 2021; Qiu et al. 2023b).

63 In the framework detailed above, the development of automatic identification techniques can  
64 significantly support the handling and processing of large raw microscopy images (Alfárez et al.  
65 2021; Faria et al. 2022). To achieve this goal, the use of Machine Learning (ML) techniques  
66 deserves attention, since these have been successfully applied in many fields of visual data

67 investigations (Petrelli and Perugini 2016; Endert et al. 2017; Acosta et al. 2019; Y. D. Wang et  
68 al. 2021; Zhou et al. 2022; Qiu et al. 2023a).

69 Among ML techniques, the developments of deep learning algorithms have drastically  
70 boosted the application of the Artificial Intelligence (AI) in many scientific fields, including image  
71 analysis and processing (Xing et al. 2018; Zhichao Liu et al. 2021). Examples are image  
72 classification (Obaid et al. 2020), object detection (Zhao et al. 2019; Wu et al. 2020) and image  
73 segmentation (Ghosh et al. 2019; Leichter et al. 2022; Tang et al. 2022). In particular, the recent  
74 development of new network algorithms in natural language processing (NLP) favored the growth  
75 of an architecture named Transformer (Vaswani et al. 2017). Transformers are at the base of the  
76 so-called foundation models (Bommasani et al. 2021), implementing the concept of “transfer  
77 learning” (Thrun and Mitchell 1995; Polat et al. 2021). The idea behind “transfer learning” is to  
78 use the “knowledge” that is learned from one task, and apply it to solve a different problem. In  
79 deep learning, the transfer learning is often achieved by the so-called “pretraining” (Bommasani  
80 et al. 2021) on large data set. More specifically, a deep learning model is typically trained to solve  
81 a non-specific task, and then adapted to the problem of interest through fine-tuning, drawn by a  
82 specific and more focused data set (Bommasani et al. 2021).

83 In this study, we investigate the application of the “transfer learning” paradigm using a  
84 Transformer known as “Swin Transformer” (Ze Liu et al. 2021). The investigated Transformer has  
85 been previously pre-trained using a large, general-purpose, public computer vision dataset. Our  
86 main aim is to tap the benefit from the “transfer learning” paradigm by fine-tuning the “Swin  
87 Transformers” in the classification of five metallic minerals (i.e., arsenopyrite, chalcopyrite, gold,  
88 pyrite and stibnite) in optical microscopy images. To achieve this goal, we set up and trained the  
89 Swin Transformer plus two widely-used Convolutional Neural Networks (CNNs) models. Then,

90 we evaluated and compared the performances of each model. Finally, we used a feature  
91 visualization technique named Class Activation Map (Zhou et al. 2016) to attempt at understanding  
92 the internal behavior of the investigated models, which is often perceived as a “black box”  
93 (Castelvecchi 2016).

## 94 **Materials**

95 Raw images were captured using optical microscopes (AXIOSCOPE-A1, Leica DM4P, and  
96 Olympus BX51) through employing cellSens Entry and Stream Essentials software under reflected  
97 light conditions. The raw data resolutions varied from 1608\*1608, 1936\*1216, and to 4800\*3600  
98 pixels. The process of collecting images involved several manual steps. Firstly, the thin-sections  
99 were observed through optical microscopy, the target minerals were located, and the mineral  
100 images were captured, particularly focusing on grains larger than 10  $\mu\text{m}$  in width to ensure accurate  
101 mineral identification under the microscope. Notably, the collected images may show different  
102 colors for the same mineral, whereas the same/similar colors are seen for different minerals. This  
103 discrepancy arises because the original data are from different deposits and the thin-sections vary  
104 in white balance and brightness.

## 105 **Data composition**

106 Microscopy images of arsenopyrite, chalcopyrite, gold, pyrite, and stibnite were selected as  
107 the research material for our target aimed at classification. The dataset consists of 481 optical  
108 microscopy images of five unprocessed metal minerals from two gold deposits in the Jiaodong  
109 Peninsula of North China (Linglong, and Longkou gold deposits) and six gold deposits in West  
110 Qinling in Central China (Jiagantan, Liba, Xiakanmucang, Yidinan, and Zaorendao gold deposits;  
111 and the Zaozigou gold stibnite deposit). A total of 159 arsenopyrite images, 128 chalcopyrite

112 images, 159 gold images, 145 pyrite images and 131 stibnite images of different sizes were used  
113 as the raw data. The characteristics of these five types of metallic minerals under microscope –  
114 that is, the information that manual classification uses—are given in Table 1.

### 115 **Dataset characteristics**

116 As stated above, metal minerals were imaged using optical microscopy. The produced images  
117 contain basic information that characterizes each phase, i.e., the reflected color, microstructural  
118 characteristics, and mineral paragenesis. These characteristics constitute the building blocks for  
119 the classification of the metal mineral phase with direct observation.

120 In detail, gold (Au; Figures 1a-d) is bright yellow under the microscope, with relatively  
121 smooth edges. Gold mainly coexists with arsenopyrite, pyrite, chalcopyrite and stibnite. Pyrite  
122 ( $\text{FeS}_2$ ; Figures 1e-h) is a homogeneous mineral appearing light yellow under the microscope with  
123 the edges of the grains being relatively smooth, as for gold. It is widely distributed with a number  
124 of minerals especially gold, arsenopyrite and chalcopyrite. In addition, the images show that  
125 chalcopyrite ( $\text{CuFeS}_2$ ; Figures 1i-l) is characterized by copper-yellow, weakly homogeneous, and  
126 often heteromorphic granular aggregates. The reflective color always shows yellow-green. Under  
127 the microscope, chalcopyrite has broken edges. Arsenopyrite ( $\text{FeAsS}$ ; Figures 1m-p) is bright  
128 white with cream, yellow, or red color (i.e., is weakly polychromatic) and radial aggregates can be  
129 observed. Arsenopyrite has smooth edges and is often arsenic-bearing pyrite and arsenopyrite.  
130 Stibnite ( $\text{Sb}_2\text{S}_3$ ; Figures 1q-t) of the white or light off-white variants can be easily confused with  
131 arsenopyrite. Similar to arsenopyrite, the strongly homogenous stibnite coexisted with pyrite with  
132 smooth edges under the microscope.

133 Gold, pyrite, and chalcopyrite are all yellow under the microscope and tend to coexist in gold  
134 deposits. The gold and pyrite grains in the dataset have relatively smooth edges, whereas  
135 chalcopyrite has broken boundaries. The arsenopyrite and stibnite all have a reflective color of  
136 gray with relatively smooth edges and similar mineral morphology.

137 From the above-mentioned features, it is obvious that large-scale manual classification by  
138 observing these metal minerals with the naked eye (or directly under the microscope) poses a  
139 significant challenge. The similarity in reflected colors and morphology complicates their  
140 classification, and would consume long of time with manual studies, especially when dealing with  
141 large volumes of images for examination, such as in batch studies.

## 142 **Methods**

### 143 **Swin Transformer and Convolutional Neural Networks**

144 **Swin Transformer.** In image analysis, “Transformers” (Dosovitskiy et al. 2020) rely on a  
145 self-attention mechanism to model the correlation between various regions within an image. The  
146 self-attention mechanism, often referred to as scaled dot-product attention, stands as a fundamental  
147 concept in the field of deep learning, allowing the model to gauge the significance of distinct  
148 elements in an input sequence, dynamically regulating their impact on the output. Compared with  
149 the local receptive field mechanism of convolution used in CNNs, transformers can learn the  
150 correlation among relatively distant areas, and capture the long-distance dependence of the whole  
151 feature map, and therefore they are characterized by a high global modeling ability (Vaswani et al.  
152 2017). The multi-head self-attention of the transformer that gives the model the ability to focus on  
153 different locations allows the model to learn relevant information in different subspaces and extract  
154 richer feature information (Devlin et al. 2019), thus alleviating the complexity of the neural

155 network. The model does not need to input all the information into the neural network for  
156 calculation, but selectively enters some task-related information into the network. However,  
157 “plain” transformers require a large amount of computation for the training (Vyas et al. 2020). The  
158 Swin Transformer algorithm aims to improve upon this by using a hierarchical approach (shifted  
159 window) to decrease the cost of computing the self-attention which is exponential in the image  
160 size. Specifically, Swin Transformer computes self-attention on a local window which is then  
161 moved across the image, and also a multi-scale feature computation method. Swin Transformer  
162 has been used as a feature extraction network in various image tasks with good results (Jiang et al.  
163 2022; Yang and Yang 2023).

164 In the present study, we optimize computational expenses by employing the Swin  
165 Transformer-base version (Ze Liu et al. 2021). The network architecture adopts a hierarchical  
166 design encompassing four stages (Figure 2). Initially, the RGB image undergoes the Patch Partition  
167 module, segmenting it into non-overlapping patches. Every pixel adjacent to 4\*4 is a Patch. The  
168 image is then flattened across the color channel dimension, reshaping the color channel values into  
169 an elongated one-dimensional vector. This sequential vector captures the color information  
170 corresponding to each pixel. Subsequently, the channel data are transformed by a Linear  
171 Embedding layer. This layer is a technique for representing images as dense embedding vectors  
172 and these vectors capture visual features of the image. Following this, feature maps of varying  
173 sizes are constructed through four stages. Except that the image first passes through a Linear  
174 Embedding layer in stage 1, for the remaining three stages, the image is subsampled through a  
175 Patch Merging layer and then through the Swin Transformer Block. The Swin Transformer-base  
176 version is described in detail by Ze Liu et al. (2021).



177       **Convolutional Neural Networks.** Convolutional neural networks (CNNs) constitute a neural  
178 network architecture that is often used to extract features from image data. The Convolutional  
179 layer is the fundamental aspect of CNNs which involves the application of convolution operation  
180 to the input data and plays a crucial role in feature extraction from images with spatial  
181 relationships. The CNNs were first proposed by LeCun et al. (1989, 2015) for handwritten digital  
182 image recognition. In the present study, two general-purpose and widely-used CNN models have  
183 been investigated, namely ResNet50 and MobileNetv2 (Sandler et al. 2018).

184       ResNet50, a 50-layer variant of ResNet (He et al. 2015) operates through five processing  
185 stages (Figure 3). Stage 1 can be considered as a preprocessing step for the input images. In detail,  
186 for a three-channel RGB input image, it performs a preliminary feature extraction via 64  
187 convolutional layers. Feature normalization is then carried out by a Batch Normalization (BN)  
188 layer that can convert interlayer outputs of a neural network into a standard format by subtracting  
189 the batch mean and then dividing by the batch's standard deviation, and effectively 'resets' the  
190 distribution of the output of the previous layer to be more efficiently processed by the subsequent  
191 layer (Chang and Chen 2015). Thus, the training convergence speed of the model is made faster  
192 and training process becomes more stable. Feature normalization is followed by a nonlinear feature  
193 mapping via the ReLU activation function layer and, next, a Maximum Pooling layer.  
194 Subsequently, the feature map size is further reduced to a quarter of the input image to reduce  
195 spatial information and parameters, increase computational speed, and reduce the risk of  
196 overfitting. In ResNet50, the remaining four stages have a similar structure, all of which are made  
197 up of different numbers of residual modules (Feng 2017). Finally, the extracted features pass  
198 through a fully connected (FC) layer to integrate features together for easy submission to the final  
199 classifier.

200 The MobileNetv2 is a common lightweight CNN characterized by the structure reported in  
201 Figure 4. The activation function ReLU is used within the MobileNetv2 because of its simplicity  
202 and efficiency of calculation. In addition, MobileNetv2 has bottleneck layers in the residual  
203 connections that obtain a representation of the input with reduced dimensionality (Sze et al. 2017).  
204 The lightweight depth-wise convolutions are used by the intermediate layer (Alain and Bengio  
205 2016) to filter features as the source of nonlinearity. MobileNetv2 has 32 filters and an initial fully  
206 connected convolution layer followed by 19 residual bottleneck layers. The fully connected  
207 convolution layer, also known as convolutional kernels, can find the most effective filters based  
208 on the task and then combine these filters into more complex patterns. In addition, the output of  
209 some neurons will be 0, which reduces the interdependence of parameters and alleviates the  
210 overfitting problem.

211 Compared with the traditional CNNs, Swin Transformer has the unique shifted window which  
212 enables the model to gain strong global modeling ability and fewer parameters (Vaswani et al.  
213 2017; Devlin et al. 2019). Due to the hierarchical approach, Swin Transformer has strong  
214 scalability in processing images of different scales. Meanwhile, the shifted window brings high  
215 computational complexity (Vyas et al. 2020). The CNNs usually have lower computational  
216 complexity and memory consumption, which can effectively extract local features in images  
217 (Zhichao Liu et al. 2021). Also, the CNNs can effectively extract local features in images through  
218 weight sharing (Abdel-Hamid et al. 2012; Miao et al. 2016).

219 The same epoch, batch size and optimizer as Swin Transformer are used in our experiments  
220 of ResNet50 and MobileNetv2.

## 221 **Classification workflow**

222 As mentioned in previous sections, we aimed at exploring the effectiveness of deep learning  
223 techniques as a possible substitute to the “by-hand” classification of metal mineral phases in  
224 images acquired by optical microscopes. To achieve our goal, we proposed a four-step workflow  
225 (Figure 5). From the raw images, we cropped the areas and selected those areas that contain only  
226 one mineral phase for further analysis (step 1); from these images, we constructed the training and  
227 test datasets, also using data augmentation techniques (step 2); then, we used the obtained datasets  
228 to train and test the investigated models (step 3); the trained models are finally used to infer the  
229 five metal mineral classes that are the object of the present study (step 4). In the following, we are  
230 going to detail the 4 steps of the proposed workflow (Figure 5).

231 Step 1: Image Stacks. We started with 481 raw light microscopy images of the five different  
232 metal mineral phases from different outcrops (see Materials Section). We next used the sliding  
233 window technique from OpenCV (Rosebrock 2015) to capture the mineral phases that are present  
234 in the image which we next cropped to equal-sized, non-overlapping, sub-images (256×256 pixels  
235 each, RGB; Figure 5). Subsequently, we removed the images that contain more than one mineral.  
236 The images spatially dominated by a single mineral were manually selected and labeled.  
237 Accordingly, 4524 images were obtained (Table 2), each being labeled with the name of the  
238 (dominant) mineral phase.

239 Step 2: Preprocessing and Data Augmentation. We resized the size of images to 224×224  
240 because of the requirement of model input. And we then divided the above-mentioned image  
241 dataset into a training set, a validation set, and a test set with a ratio of 3:1:1 (see Table 2). The  
242 scope of the training set is to “educate” the model by determining the weight and bias learning  
243 parameters. The validation set is used to tune hyperparameters and check whether the effect of

244 model training goes in a “good” or “bad” direction. In addition, the validation set is used and data  
245 is unseen during the training process. Finally, the test set is used to evaluate the generalization  
246 ability of the final model. In addition, the samples in the test set were selected from different ore  
247 deposits than those in the training set, which alleviates the potential issue of data leakage.

248       However, due to the limited number of images in the dataset, training can be challenging. To  
249 improve the model’s generalization ability and reduce overfitting, we increased the size and the  
250 robustness of the training set and introduced variations that could be found in “real world data”  
251 using six typical offline data augmentation methods (Supplemental Materials). The first is named  
252 “random erasing” (Zhong et al. 2020). It consists of randomly selecting a rectangular region in an  
253 image and replacing its pixels with random values. This procedure generates new training images  
254 with various levels of occlusion, which, when used for training, reduces the risk of over-fitting and  
255 also makes the model less sensitive to occlusions (i.e., missing portions). The second approach is  
256 the “flipping”. It consists of mirroring the images both horizontally and vertically, along the  
257 vertical and horizontal axes, respectively. The third augmentation method is named “brightness  
258 adjust”. To note, the coloring of a picture can be set using three parameters: hue (H), saturation  
259 (S) and value (V), with the latter mainly governing the brightness. By using the Auto Gamma  
260 Correction method, a non-linear operation  $S = T(R) = R^\gamma$  (where S and R are the values of  
261 brightness in output and original image, respectively, that are mapped to [0 1]) is to automatically  
262 lighten or darken the image (Babakhani and Zarei 2015). The fourth approach is “random zoom”,  
263 which zooms into an image at a random location within the image. The fifth is “random contrast”,  
264 which adjusts the contrast of the images by a random factor. And the last is “random saturation”,  
265 which can adjust the saturation of images by a random factor. These methods can also improve the

266 model's ability to classify based on the color (Supplemental Materials). At the end of the  
267 augmentation process, we increased the number of the training set to 18991 images.

268       Step 3: Model Training. For the present study, we trained two standard CNNs (i.e., ResNet50  
269 and MobileNetv2) and a "Swin Transformer" algorithm (Figure 5). These architectures are  
270 followed by the max pooling and a fully connected (FC) layer. From the latter, a softmax function  
271 performs the final prediction, selecting the category with the largest softmax value. We trained the  
272 three models using adam gradient descent algorithm (Kingma and Ba 2017), using the first-order  
273 momentum parameter of 0.9, the second-order momentum parameter of 0.999, a learning rate of  
274 1e-6 and a batch size of 16, respectively. The task was set for 20 epochs, each of which contains  
275 136 iterations. At each training iteration, the image fed into the model is forward-propagated to  
276 the output layer, after which the difference between the ground-truth label and predictive label is  
277 measured by a Cross Entropy (CE) loss function. The loss value is then reduced by back-  
278 propagating and updating the model's parameters. To accelerate the training convergence and  
279 possibly increase the generalization capabilities of the models, we used a "transfer learning"  
280 approach by initializing the weights of the models to those from ImageNet-1k (Deng et al. 2009).  
281 In the process of model evaluation, the accuracy for the validation set is calculated after every  
282 epoch and the model's final accuracy is the highest among 20 epochs (Wu and Chen 2015; Ruby  
283 and Yendapalli 2020; Zhong et al. 2020).

284       Step 4: Getting the Output. We finally obtained a metal mineral classifier and we evaluated  
285 its performance.

286 **Model evaluation**

287 For the evaluation of the investigated models, we utilized 5 metrics, i.e., accuracy, precision,  
288 recall, F1-score, and training loss, defined as follows.

289 The model accuracy is defined as

290 
$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \quad (1)$$

291 TP, TN, FP, and FN are the true positive, true negative, false positive, and false negative  
292 occurrences, respectively. The accuracy is the most common evaluation index in deep learning for  
293 classifiers, i.e., deep learning models used in classification tasks, due to its simple intuition.

294 The model precision is defined as

295 
$$Precision = \frac{TP}{TP + FP} \quad (2)$$

296 The precision is the proportion of positive samples predicted correctly by the model from all  
297 samples predicted as positive.

298 The model recall is defined as

299 
$$Recall = \frac{TP}{TP + FN} \quad (3)$$

300 The recall is a coverage measure, representing the classification accuracy of positive samples.

301 The F1-score is given by

302 
$$F1 = 2 \times \left[ \frac{Precision \times Recall}{Precision + Recall} \right] \quad (4)$$

303 The F1-score is the harmonic mean between precision and recall. A macro F1-score computes  
304 the metric for each category independently and then takes the average (all categories are treated  
305 equally).

306 Finally, the Cross Entropy function is used to calculate the training loss of the model, which  
307 is defined as

$$308 \quad CE(p, q) = -\sum_{i=1}^C p_i \log(q_i) \quad (5)$$

309 where  $C$  denotes the total number of classes,  $p_i$  denotes the  $i$ -th prediction class probability,  
310 and  $q_i$  denotes the  $i$ -th true class of training samples. The smaller the Cross Entropy (CE) loss is,  
311 the distributions of the two probabilities are approximately close, indicating that the model has a  
312 good performance.

### 313 **Transfer learning**

314 To get benefit from the “transfer learning” paradigm (Torrey and Shavlik 2010; Zhichao Liu  
315 et al. 2021), we set the initial weights of the investigated models to the pre-trained values that have  
316 been obtained after a full training on the natural ImageNet-1k dataset. The ImageNet-1k is a large  
317 public computer vision dataset that is often used as a benchmark to evaluate the performances of  
318 different deep learning models. It consists of 10 million images, characterized by thousand  
319 categories (Deng et al. 2009). The weights trained by ImageNet-1k, when used as the pre-training  
320 weight initialization, can quickly extract the shallow general image features (such as shape,  
321 brightness, and size of underlying image structures), thus future improving the initial accuracy of  
322 the model and accelerating the convergence of training models (Deng et al. 2009; Torrey and  
323 Shavlik 2010). Pre-trained weights for the ResNet50, MobileNetv2 and Swin Transformer are  
324 publicly available at the following repository: <https://download.pytorch.org/models>.

## 325 **Class Activation Map**

326 Although deep learning models are often characterized by good performances, they are  
327 subject to criticism because of their “black box” nature. To unblur the “black box” nature of the  
328 proposed models, we adopted the Class Activation Map (CAM; Zhou et al. 2016), also known as  
329 Class Thermal Maps. The CAM is a feature visualization technique that aims at highlighting the  
330 contribution of the different image regions to a given classification outcome. In detail, a CAM for  
331 a particular class of objects highlights the image regions used by the model to identify the specific  
332 class and shows which feature maps the model is based on for classification. Using a network  
333 architecture comprising convolutional layers, the feature map is extracted and the feather map up-  
334 sampled, which could be used as mask information to obtain the model's response value to the  
335 image in the target class. By linearly weighting the feature map with the obtained response values,  
336 the visual CAM mappings are obtained (Wang et al. 2020). In our specific case, the areas  
337 characterized by low- and high-discriminative powers have been highlighted by progressively  
338 shifting the thermal map from the blue to the red color. CAM is fully described in Zhou (2016).

## 339 **Codes and Libraries**

340 All the tasks are implemented using Python (<https://www.python.org>; version 3.8) and  
341 Pytorch (<https://pytorch.org>; version 1.8.1), and finished by Intel(R) Xeon(R) CPU E5-2630 v4 @  
342 2.20GHz and NVIDIA GeForce RTX 3090 GPU. All the codes and the dataset are available at the  
343 <https://doi.org/10.5281/zenodo.10441351> repository.

344 The following libraries were used to complete the code: Pytorch (<https://pytorch.org/>) for  
345 computing tensors on graphics processing units; NumPy (<https://numpy.org/>) for data analysis;  
346 TorchAudio (<https://pytorch.org/audio/stable/index.html>) and SciPy (<https://scipy.org/>) for data  
347 processing functions; TensorBoard (<https://tensorflow.google.cn/tensorboard>) for data



348 visualization; OpenCV (<https://opencv.org/>) and Pillow (<https://pypi.org/project/Pillow/>) for  
349 image processing; Torchvision (<https://pytorch.org/vision/stable/index.html>) for image  
350 classification; timm (<https://timm.fast.ai/>) for loading image model; Safetensors  
351 (<https://pypi.org/project/safetensors/>) for parameters and weights saving; tqdm  
352 (<https://pypi.org/project/tqdm/>) for progress prompt; Matplotlib (<https://matplotlib.org/>) for  
353 plotting the diagrams.

## 354 **Results**

355 Figure 6 displays the evolution of the training loss and validation accuracy for the Swin  
356 Transformer, ResNet50, and MobileNetv2, respectively.

### 357 **Swin Transformer**

358 For the Swin Transformer, the accuracy of the validation set gradually increases from 0.74 at  
359 the beginning of the training to 0.92 after 16 epochs (Figure 6a). After seven epochs, the accuracy  
360 reached 0.90. On the training set, the Cross Entropy (CE) loss which is a function that minimizes  
361 the model's loss is initially 1.42 (Figure 6b). The CE loss on the training set, initially at 1.60,  
362 decreases considerably during the seven epochs and slowly decreases over the course of the  
363 following epochs, reaching a minimum value of 0.01 in twentieth epoch (Figure 6b).

364 The average F1-score of each category for the Swin Transformer model is 0.92 (Table 3).  
365 Further analysis of the performance of each category shows that gold has the best classification  
366 performance with the F1-score of 0.98, the recall of 1.00 and the precision of 0.96 while pyrite has  
367 the lowest classification performance with the recall of 0.81. Figure 7a further details these insights  
368 by showing the confusion matrix for this model.

369 **Convolutional Neural Networks**

370 **ResNet50.** The accuracy of ResNet50 on the validation set gradually increases from 0.40 to  
371 0.91, after 20 epochs (Figure 6a). After about 13 epochs, the accuracy of the validation set has  
372 reached 0.90, and then it begins to stabilize gradually. Initially, the CE loss on the training set is  
373 1.58. During seven epochs, the training loss decreases greatly and drops below 0.1. The training  
374 loss decreases slowly in later epochs and reaches the minimum value of 0.01 after twenty epochs  
375 (Figure 6b).

376 The average F1-score of each category for the ResNet50 is 0.90 (Table 3). Further analysis  
377 of the performance parameters for each class shows that stibnite has the best classification  
378 performance whose F1-score is 0.98, the recall is 0.98 and the precision is 0.97, whereas  
379 chalcopyrite and pyrite have the worst performance with an F1-scores of 0.84 and 0.80,  
380 respectively. As can be found from the confusion matrix of ResNet50 (Figure 7b), chalcopyrite  
381 and pyrite are not predicted well, while gold is the best.

382 **MobileNetv2.** For the MobileNetv2, the accuracy of the validation set gradually increases  
383 from 0.34 to 0.84, achieved after 20 epochs (Figure 6a). After six epochs, the accuracy of the  
384 validation set has reached 0.81, after which it begins to stabilize gradually. The CE loss on the  
385 training set, initially at 1.60, continues to reduce in later epochs until it reaches a minimum value  
386 of 0.01. (Figure 6b).

387 The average F1-score for each class in the MobileNetv2 is 0.81 (Table 3). Further analysis of  
388 the performance of each class shows that gold is the best classified mineral, while pyrite and  
389 stibnite are the worst, with recall of 0.61, and 0.62, respectively. The confusion matrix (Figure 7c)  
390 shows that the MobileNetv2 can easily predict arsenopyrite and gold.

391

## Discussion

### 392 **Model classification performance**

393 The results of the present study demonstrate that the Swin Transformer is characterized by an  
394 excellent prediction performance and a higher accuracy than the other tested models (Table 3).  
395 Comparing the specific scores of the three models, the Swin Transformer greatly improves the F1-  
396 scores of chalcopyrite, gold and pyrite, which are difficult to classify by the investigated CNNs  
397 (Table 3). This occurrence results in final average class accuracies of 0.92, 0.91, and 0.81 for the  
398 Swin Transformer, ResNet50, and MobileNetv2, respectively. The Swin Transformer also  
399 provides the lowest final training loss (Figure 6).

400 Moreover, misclassification occurrences of chalcopyrite and pyrite, often recognized as gold  
401 by ResNet50 (F1-score equal to 0.84 and 0.80, respectively; Table 3), MobileNetv2 (F1-score  
402 equal to 0.83 and 0.75, respectively; Table 3) as highlighted in Figure 7, were greatly reduced by  
403 the use of the Swin Transformer (F1-score of chalcopyrite and pyrite exceeding 0.85). It is also  
404 seen that stibnite is much less likely to be misclassified as arsenopyrite by the Swin Transformer  
405 than the MobileNetv2 (Figure 7). As a drawback, the confusion matrix demonstrates that stibnite  
406 is more likely misclassified as chalcopyrite by the Swin Transformer than the Resnet50 (Figure 7).  
407 This occurrence results in a slightly lower precision value of chalcopyrite of Swin Transformer  
408 (i.e., 0.85) than those characterizing the ResNet50 (i.e., 0.97). By analyzing the cause of the  
409 classification error, it can be inferred that the model easily classifies pyrites as chalcopyrite,  
410 probably due to the limited quality of the input images, thus the two minerals in the images have  
411 a similar yellow color (Figures 8a, b). This feature can confuse the models. Furthermore,  
412 arsenopyrite and stibnite are often misclassified by all models. The visual analysis of these  
413 situations (Figures 8c, d) shows that there are two main reasons for misclassification: (1) As

414 stibnite images are from different samples with different image collection parameters, and they  
415 show a variety of colors some of which are similar to the grey reflection color of the arsenopyrite.  
416 (2) For some samples, both minerals have similar crystal forms from euhedral to subhedral, which  
417 can confuse the feature identification of the network. Simultaneously, pyrite and arsenopyrite also  
418 have similar reflection colors which are not easy to distinguish (Figures 8e, f). Also, other minerals  
419 present in the image cause interference, leading to model classification error (Figures 8e, f).  
420 However, these similarities do not affect the overall classification capabilities of the Swin  
421 Transformer, which clearly outperforms the investigated CNNs in most of the scores for the single  
422 classes and all the average performance metrics (Table 3 and Figure 7). In conclusion, our study  
423 supports the Swin Transformer as a metal mineral classifier (abbreviated ST-MMC).

#### 424 **Transfer learning in optical microscopy for the study of metal minerals**

425 A large number of studies have shown that using the transfer learning paradigm to set the  
426 initial weights of a model before starting the training can effectively help in achieving the  
427 convergence. Transfer learning can fine-tune the parameters of the entire model to get initial high  
428 accuracy and a low loss value. (Figure 9; Supplemental Materials). Also, it allows for improving  
429 the generalization capability of a model (Kora Venu 2022). As an example, a number of studies  
430 demonstrated that using the “knowledge” acquired on natural images effectively improves the  
431 capabilities on a model in solving specific problems like the processing of medical or remote  
432 sensing images, even when using limited training sets (Xie et al. 2016; Raghu et al. 2019; Kora et  
433 al. 2021). Collecting metal mineral images with optical microscopes is a time-consuming task,  
434 thus resulting in a limited dataset size. As reported in the method section, we adopted the “pre-  
435 trained” weights for the investigated models deriving from a training on the ImageNet-1k dataset.  
436 As highlighted in Figure 6, the accuracy of the first validation, i.e., deriving from the pretraining

437 only, was 0.34, 0.40, and 0.74 for the MobileNetv2, ResNet50 and the Swin Transformer,  
438 respectively. The Swin Transformer has a greater response to transfer learning and a higher initial  
439 accuracy.

440 To further outline the added value of transfer learning in achieving a solution for the problem  
441 investigated in the present manuscript, we trained the Swin Transformer without pre-trained  
442 weights. The training epoch, batch size, and other parameters of both models were the same. As a  
443 result, with the same number of training iterations, we obtained an initial and maximum accuracy  
444 of 0.56 and 0.88, respectively (Figure 9a), less than the accuracy achieved with the support of the  
445 “transfer learning” paradigm, i.e., 0.92. And the use of “transfer learning” paradigm supports lower  
446 initial and minimum loss values, which are 1.42 and 0.01, respectively (Figure 9b).

#### 447 **Model interpretation**

448 Figure 10 shows five images that have been fed as unknowns, classified by the three models  
449 investigated in the present manuscript and output the probabilities. For each image, a blue-to-red  
450 heatmap points to the contribution of the different regions to the classification output. In detail,  
451 Figure 10 highlights that, for minerals with broken edges such as chalcopyrite (Figure 10a),  
452 thermal maps specifically focus on mineral edges. For the other cases (Figures 10b, c, d), the  
453 thermal maps highlight different regions, often focusing on the edges.

454 Based on the evidence reported above, it can be inferred that the shape of edges effectively  
455 influences the classification and the networks pay attention to their smoothness or sharpness.  
456 Moreover, the occurrences of misclassifications can now be better explained: Arsenopyrite and  
457 stibnite are both grey in color, and they also share smooth edges (Figures 8c, d). Also, both stibnite  
458 and arsenopyrite have void development on their surfaces. Despite the insights provided by CAMs

459 do not directly lead us to improved models, they point to the causes that generate misclassifications  
460 and, therefore, suggest a direction for the possible improvement.

461 Notably, the different investigated networks show significant differences in their CAMs for  
462 the same inputs (Figure 10). The thermal maps of ResNet50 and MobileNetv2 almost focus on the  
463 edge of the minerals, and not inside them. As a consequence, the reflected color and texture may  
464 not be the most important distinctive features of these two models. However, CAMs for the Swin  
465 Transformer also cover the interior of the mineral, rather than just the edge, which suggests that  
466 ST-MMC effectively uses the reflected color and texture of the minerals for its inference, and it  
467 has better global performance. In the thermal map, the mineral area of middle and edge of ST-  
468 MMC are redder than the other two CNNs, indicating that these domains have stronger model  
469 response, which reveals that mineral reflection color and texture contribute more to the  
470 classification output of the model. The Swin Transformer also achieved a classification predict  
471 probability over 0.95 for unknown minerals, significantly outperforming the other two CNNs  
472 (Figure 10). This effectiveness in handling unknown samples also demonstrates its capability for  
473 efficient batch image processing.

## 474 **Implications**

475 Large-dimensional image analyses are dominantly based on digital image datasets, the  
476 automatic identification of the optical microscopic data is still poorly examined, and the mineral  
477 image data is also difficult to collect. Deep learning-based approach (Swin Transformer) with the  
478 transfer learning paradigm fully explores the information of different metal mineral phases to  
479 produce a well-behaved mineral classifier with high accuracy and strong global ability. To  
480 circumvent the ‘black box’ problem commonly associated with deep learning models, CAM (Class  
481 Activation Map) tool was introduced to explain individual predictions. With the increasing amount

482 of high-throughput mineral image data produced by modern analytical techniques, our ST-MMC  
483 offers the potential to make more data driven decisions such as transparent minerals classification.  
484 Moreover, the “transfer learning” paradigm on large images captured by optical microscopy, will  
485 possibly liberate researchers from tiresome labor, sharpen the accuracy, and increase the  
486 productivity. More widely, the use of “transfer learning” may disclose new perspectives in  
487 petrology and mineralogy, possibly providing a paradigm shift over the current applications of  
488 deep learning in petrology and mineralogy.

#### 489 **Acknowledgments**

490 We gratefully acknowledge editors for handling our manuscript, and the constructive  
491 comments of the anonymous reviewers. This research was financially supported by the National  
492 Key Research and Development Program (SQ2023YFE0103286), the National Natural Science  
493 Foundation (42261134535 and 42072087), the Frontiers Science Center for Deep-time Digital  
494 Earth (2652023001), the 111 Project of the Ministry of Science and Technology (BP0719021) and  
495 the Shandong Provincial Engineering Laboratory of Application and Development of Big Data for  
496 Deep Gold Exploration (SDK202211 and SDK202214). MP kindly acknowledge the MUR  
497 PRIN2020 project “Dynamics and timescales of volcanic plumbing systems: a multidisciplinary  
498 approach to a multifaceted problem” (202037YPCZ\_001).

#### 499 **References**

500 Abdel-Hamid, O., Mohamed, A., Jiang, H., and Penn, G. (2012) Applying Convolutional Neural  
501 Networks concepts to hybrid NN-HMM model for speech recognition. In 2012 IEEE  
502 International Conference on Acoustics, Speech and Signal Processing (ICASSP) pp. 4277–

- 503 4280. Presented at the 2012 IEEE International Conference on Acoustics, Speech and Signal  
504 Processing (ICASSP), <https://doi.org/10.1109/ICASSP.2012.6288864>.
- 505 Acosta, I.C.C., Khodadadzadeh, M., Tusa, L., Ghamisi, P., and Gloaguen, R. (2019) A Machine  
506 Learning Framework for Drill-Core Mineral Mapping Using Hyperspectral and High-  
507 Resolution Mineralogical Data Fusion. IEEE Journal of Selected Topics in Applied Earth  
508 Observations and Remote Sensing, 12, 4829–4842, [https://doi.org/10.1007/s00521-021-](https://doi.org/10.1007/s00521-021-05849-3)  
509 05849-3.
- 510 Alain, G., and Bengio, Y. (2016) Understanding intermediate layers using linear classifier probes.  
511 arXiv preprint arXiv:1610.01644, <https://doi.org/10.48550/arXiv.1610.01644>.
- 512 Alférez, G.H., Vázquez, E.L., Martínez Ardila, A.M., and Clausen, B.L. (2021) Automatic  
513 classification of plutonic rocks with deep learning. Applied Computing and Geosciences, 10,  
514 100061, <https://doi.org/10.1016/j.acags.2021.100061>.
- 515 Azeuda Ndonfack, K.I., Xie, Y., and Goldfarb, R. (2022) Gold occurrences of the Woumbou–  
516 Colomine–Kette district, eastern Cameroon: ore-forming constraints from petrography,  
517 SEM–CL imagery, fluid inclusions, and C–O–H–S isotopes. Mineralium Deposita, 57, 83–  
518 105, <https://doi.org/10.1007/s00126-021-01050-7>.
- 519 Babakhani, P., and Zarei, P. (2015) Automatic gamma correction based on average of brightness,  
520 4, 156-159.
- 521 Bommasani, R., Hudson, D.A., Adeli, E., Altman, R., Arora, S., von Arx, S., Bernstein, M.S.,  
522 Bohg, J., Bosselut, A., Brunskill, E., and others (2021) On the Opportunities and Risks of  
523 Foundation Models. arXiv preprint arXiv:2108.07258,  
524 <https://doi.org/10.48550/arXiv.2108.07258>.
- 525 Cameron, E.N. (1961) Ore microscopy.



- 526 Castelveccchi, D. (2016) Can we open the black box of AI? *Nature News*, 538, 20.
- 527 Chang, J.-R., and Chen, Y.-S. (2015) Batch-normalized maxout network in network. arXiv  
528 preprint arXiv:1511.02583, <https://doi.org/10.48550/arXiv.1511.02583>.
- 529 Chen, Z., Chen, D., and Li, C. (1979) Color index of the reflection color of ore minerals. *Acta*  
530 *Petrologica Sinica*, 219–233, <https://doi.org/10.19762/j.cnki.dizhixuebao.1979.03.005>.
- 531 Craig, J.R., Vaughan, D.J., and Hagni, R.D. (1981) *Ore microscopy and ore petrography* Vol. 406.  
532 Wiley New York.
- 533 Criddle, A.J., and Stanley, C.J. (2012) *Quantitative data file for ore minerals*. Springer Science &  
534 Business Media.
- 535 Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., and Fei-Fei, L. (2009) ImageNet: A large-scale  
536 hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern*  
537 *Recognition* pp. 248–255. Presented at the 2009 IEEE Conference on Computer Vision and  
538 *Pattern Recognition*, <https://doi.org/10.1109/CVPR.2009.5206848>.
- 539 Deng, J., Yang, L.-Q., Groves, D.I., Zhang, L., Qiu, K.-F., and Wang, Q.-F. (2020a) An integrated  
540 mineral system model for the gold deposits of the giant Jiaodong province, eastern China.  
541 *Earth-Science Reviews*, 208, 103274, <https://doi.org/10.1016/j.earscirev.2020.103274>.
- 542 Deng, J., Qiu, K.-F., Wang, Q.-F., Goldfarb, R., Yang, L.-Q., Zi, J.-W., Geng, J.-Z., and Ma, Y.  
543 (2020b) In situ dating of hydrothermal monazite and implications for the geodynamic  
544 controls on ore formation in the Jiaodong gold province, eastern China. *Economic Geology*,  
545 115, 671–685, <https://doi.org/10.5382/econgeo.4711>.
- 546 Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. (2018) BERT: Pre-training of Deep  
547 Bidirectional Transformers for Language Understanding. arXiv preprint arXiv:1810.04805,  
548 <https://doi.org/10.48550/arXiv.1810.04805>.

- 549 dos Anjos, C.E., Avila, M.R., Vasconcelos, A.G., Pereira Neta, A.M., Medeiros, L.C., Evsukoff,  
550 A.G., Surmas, R., and Landau, L. (2021) Deep learning for lithological classification of  
551 carbonate rock micro-CT images. *Computational Geosciences*, 25, 971–983,  
552 <https://doi.org/10.1007/s10596-021-10033-6>.
- 553 Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani,  
554 M., Minderer, M., Heigold, G., Gelly, S., and others (2020) An Image is Worth 16x16 Words:  
555 Transformers for Image Recognition at Scale. arXiv preprint arXiv:2010.11929,  
556 <https://doi.org/10.48550/arXiv.2010.11929>.
- 557 Endert, A., Ribarsky, W., Turkay, C., Wong, B.L.W., Nabney, I., Blanco, I.D., and Rossi, F. (2017)  
558 The State of the Art in Integrating Machine Learning into Visual Analytics. *Computer*  
559 *Graphics Forum*, 36, 458–486, <https://doi.org/10.1111/cgf.13092>.
- 560 Faria, E.L., Coelho, Juliana.M., Matos, T.F., Santos, B.C.C., Trevizan, W.A., Gonzalez, J.L., Bom,  
561 C.R., de Albuquerque, Márcio P., and de Albuquerque, Marcelo P. (2022) Lithology  
562 identification in carbonate thin section images of the Brazilian pre-salt reservoirs by the  
563 computational vision and deep learning. *Computational Geosciences*, 26, 1537–1547,  
564 <https://doi.org/10.1007/s10596-022-10168-0>.
- 565 Feng, V. (2017) An overview of resnet and its variants. *Towards data science*, 2.
- 566 Fu, Y., and Aldrich, C. (2019) Quantitative Ore Texture Analysis with Convolutional Neural  
567 Networks. *IFAC-PapersOnLine*, 52, 99–104, <https://doi.org/10.1016/j.ifacol.2019.09.171>.
- 568 Ghosh, S., Das, N., Das, I., and Maulik, U. (2019) Understanding Deep Learning Techniques for  
569 Image Segmentation. *ACM Computing Surveys*, 52, 73:1-73:35,  
570 <https://doi.org/10.1145/3329784>.

- 571 He, K., Zhang, X., Ren, S., and Sun, J. (2015) Deep Residual Learning for Image Recognition.  
572 arXiv preprint arXiv:1512.03385, <https://doi.org/10.48550/arXiv.1512.03385>.
- 573 Jiang, Y., Zhang, Y., Lin, X., Dong, J., Cheng, T., and Liang, J. (2022) SwinBTS: A Method for  
574 3D Multimodal Brain Tumor Segmentation Using Swin Transformer. *Brain Sciences*, 12,  
575 797, <https://doi.org/10.3390/brainsci12060797>.
- 576 Jungmann, M., Pape, H., Wißkirchen, P., Clauser, C., and Berlage, T. (2014) Segmentation of thin  
577 section images for grain size analysis using region competition and edge-weighted region  
578 merging. *Computers & Geosciences*, 72, 33–48,  
579 <https://doi.org/10.1016/j.cageo.2014.07.002>.
- 580 Kingma, D.P., and Ba, J. (2017, January 29) Adam: A Method for Stochastic Optimization. arXiv.
- 581 Kora, P., Ooi, C.P., Faust, O., Raghavendra, U., Gudigar, A., Chan, W.Y., Meenakshi, K., Swaraja,  
582 K., Plawiak, P., and Rajendra Acharya, U. (2021) Transfer learning techniques for medical  
583 image analysis: A review. *Biocybernetics and Biomedical Engineering*, 42, 79–107,  
584 <https://doi.org/10.1016/j.bbe.2021.11.004>.
- 585 Kora Venu, S. (2022) Improving the Generalization of Deep Learning Classification Models in  
586 Medical Imaging Using Transfer Learning and Generative Adversarial Networks. In A.P.  
587 Rocha, L. Steels, and J. van den Herik, Eds., *Agents and Artificial Intelligence* pp. 218–235.  
588 Springer International Publishing, Cham, [https://doi.org/10.1007/978-3-031-10161-8\\_12](https://doi.org/10.1007/978-3-031-10161-8_12).
- 589 LeCun, Y., Boser, B., Denker, J.S., Henderson, D., Howard, R.E., Hubbard, W., and Jackel, L.D.  
590 (1989) Backpropagation Applied to Handwritten Zip Code Recognition. *Neural*  
591 *Computation*, 1, 541–551. Presented at the Neural Computation,  
592 <https://doi.org/10.1162/neco.1989.1.4.541>.

- 593 LeCun, Y., Bengio, Y., and Hinton, G. (2015) Deep learning. *Nature*, 521, 436–444,  
594 <https://doi.org/10.1038/nature14539>.
- 595 Leichter, A., Almeev, R.R., Wittich, D., Beckmann, P., Rottensteiner, F., Holtz, F., and Sester, M.  
596 (2022) Automated Segmentation of Olivine Phenocrysts in a Volcanic Rock Thin Section  
597 Using a Fully Convolutional Neural Network. *Frontiers in Earth Science*, 10, 740638,  
598 <https://doi.org/10.3389/feart.2022.740638>.
- 599 Liu, Zhichao, Jin, L., Chen, J., Fang, Q., Ablameyko, S., Yin, Z., and Xu, Y. (2021) A survey on  
600 applications of deep learning in microscopy image analysis. *Computers in Biology and*  
601 *Medicine*, 134, 104523, <https://doi.org/10.1016/j.compbimed.2021.104523>.
- 602 Liu, Ze, Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., and Guo, B. (2021) Swin  
603 Transformer: Hierarchical Vision Transformer Using Shifted Windows pp. 10012–10022.  
604 Presented at the Proceedings of the IEEE/CVF International Conference on Computer Vision,  
605 <https://doi.org/10.48550/arXiv.2103.14030>.
- 606 Lu, J., and Peng, X. (2010) Manual of microscopic identification of metal minerals. Geology Press,  
607 Beijing.
- 608 Miao, S., Wang, Z.J., and Liao, R. (2016) A CNN Regression Approach for Real-Time 2D/3D  
609 Registration. *IEEE Transactions on Medical Imaging*, 35, 1352–1363. Presented at the IEEE  
610 *Transactions on Medical Imaging*, <https://doi.org/10.1109/TMI.2016.2521800>.
- 611 Młynarczuk, M., Górszczyk, A., and Ślipek, B. (2013) The application of pattern recognition in  
612 the automatic classification of microscopic rock images. *Computers & Geosciences*, 60, 126–  
613 133, <https://doi.org/10.1016/j.cageo.2013.07.015>.

- 614 Obaid, K.B., Zeebaree, S.R.M., and Ahmed, O.M. (2020) Deep Learning Models Based on Image  
615 Classification: A Review. *International Journal of Science and Business*, 4, 75–81,  
616 <https://doi.org/10.5281/zenodo.4108433>.
- 617 Petrelli, M. (2021) *Introduction to Python in Earth Science Data Analysis: From Descriptive*  
618 *Statistics to Machine Learning*. Springer Nature.
- 619 Petrelli, M., and Perugini, D. (2016) Solving petrological problems through machine learning: the  
620 study case of tectonic discrimination using geochemical and isotopic data. *Contributions to*  
621 *Mineralogy and Petrology*, 171, 81, <https://doi.org/10.1007/s00410-016-1292-2>.
- 622 Picot, P., and Johan, Z. (1977) *Atlas des Mineraux Metalliques*.
- 623 Piller, H. (1966) Colour measurements in ore-microscopy. *Mineralium Deposita*, 1, 175–192,  
624 <https://doi.org/10.1007/BF00204546>.
- 625 ——— (2012) *Microscope photometry*. Springer Science & Business Media.
- 626 Polat, Ö., Polat, A., and Ekici, T. (2021) Automatic classification of volcanic rocks from thin  
627 section images using transfer learning networks. *Neural Computing and Applications*, 33,  
628 11531–11540, <https://doi.org/10.1007/s00521-021-05849-3>.
- 629 Qiu, K., Zhou, T., Chew, D., Hou, Z., Müller, A., Yu, H., Lee, R.G., Chen, H., and Deng, J. (2023)  
630 Apatite trace element composition as an indicator of ore deposit types: a machine learning  
631 approach. *American Mineralogist*, <http://doi.org/10.2138/am-2022-8805>.
- 632 Qiu, K.-F., Deng, J., Laflamme, C., Long, Z.-Y., Wan, R.-Q., Moynier, F., Yu, H.-C., Zhang, J.-  
633 Y., Ding, Z.-J., and Goldfarb, R. (2023a) Giant Mesozoic gold ores derived from subducted  
634 oceanic slab and overlying sediments. *Geochimica et Cosmochimica Acta*, 343, 133–141,  
635 <https://doi.org/10.1016/j.gca.2023.01.002>.

- 636 Qiu, K.-F., Deng, J., Sai, S.-X., Yu, H.-C., Tamer, M.T., Ding, Z.-J., Yu, X.-F., and Goldfarb, R.  
637 (2023b) Low-Temperature Thermochronology for Defining the Tectonic Controls on  
638 Heterogeneous Gold Endowment Across the Jiaodong Peninsula, Eastern China. *Tectonics*,  
639 42, e2022TC007669, <https://doi.org/10.1029/2022TC007669>.
- 640 Raghu, M., Zhang, C., Kleinberg, J., and Bengio, S. (2019) Transfusion: Understanding Transfer  
641 Learning for Medical Imaging. In *Advances in Neural Information Processing Systems Vol.*  
642 32. Curran Associates, Inc, <https://doi.org/10.48550/arXiv.1902.07208>.
- 643 Ramdohr, P. (2013) *The ore minerals and their intergrowths*. Elsevier.
- 644 Rosebrock, A. (2015) Sliding windows for object detection with python and opencv.  
645 PylmageSearch, Navigation.
- 646 Ruby, U., and Yendapalli, V. (2020) Binary cross entropy with deep learning technique for Image  
647 classification. *International Journal of Advanced Trends in Computer Science and*  
648 *Engineering*, 9.
- 649 Sandler, M., Howard, A., Zhu, M., Zhmoginov, A., and Chen, L.-C. (2018) MobileNetV2: Inverted  
650 Residuals and Linear Bottlenecks pp. 4510–4520. Presented at the Proceedings of the IEEE  
651 Conference on Computer Vision and Pattern  
652 Recognition, <https://doi.org/10.48550/arXiv.1801.04381>.
- 653 Santosh, M. (2010) Assembling North China Craton within the Columbia supercontinent: The role  
654 of double-sided subduction. *Precambrian Research*, 178, 149–167,  
655 <https://doi.org/10.1016/j.precamres.2010.02.003>.
- 656 Santosh, M., Wilde, S.A., and Li, J.H. (2007) Timing of Paleoproterozoic ultrahigh-temperature  
657 metamorphism in the North China Craton: Evidence from SHRIMP U–Pb zircon

- 658 geochronology. Precambrian Research, 159, 178–196,  
659 <https://doi.org/10.1016/j.precamres.2007.06.006>.
- 660 Santosh, M., Maruyama, S., and Sato, K. (2009) Anatomy of a Cambrian suture in Gondwana:  
661 Pacific-type orogeny in southern India? Gondwana Research, 16, 321–341,  
662 <https://doi.org/10.1016/j.gr.2008.12.012>.
- 663 Schneider, C.A., Rasband, W.S., and Eliceiri, K.W. (2012) NIH Image to ImageJ: 25 years of  
664 image analysis. Nature methods, 9, 671–675, <https://doi.org/10.1038/nmeth.2089>.
- 665 Schrader, D.L., and Zega, T.J. (2019) Petrographic and compositional indicators of formation and  
666 alteration conditions from LL chondrite sulfides. Geochimica et Cosmochimica Acta, 264,  
667 165–179, <https://doi.org/10.1016/j.gca.2019.08.015>.
- 668 Shang, J., and Lin, J. (1990) Gold minerals and their occurrence. Journal of Changchun University  
669 of Earth Science, 20, 273–278.
- 670 Sheldrake, T., and Higgins, O. (2021) Classification, segmentation and correlation of zoned  
671 minerals. Computers & Geosciences, 156, 104876,  
672 <https://doi.org/10.1016/j.cageo.2021.104876>.
- 673 Su, C., Xu, S., Zhu, K., and Zhang, X. (2020) Rock classification in petrographic thin section  
674 images based on concatenated convolutional neural networks. Earth Science Informatics, 13,  
675 1477–1484, <https://doi.org/10.1007/s12145-020-00505-1>.
- 676 Sze, V., Chen, Y.-H., Yang, T.-J., and Emer, J.S. (2017) Efficient Processing of Deep Neural  
677 Networks: A Tutorial and Survey. Proceedings of the IEEE, 105, 2295–2329. Presented at  
678 the Proceedings of the IEEE, <https://doi.org/10.1109/jproc.2017.2761740>.

- 679 Tang, K., Wang, Y.D., Mostaghimi, P., Knackstedt, M., Hargrave, C., and Armstrong, R.T. (2022)  
680 Deep convolutional neural network for 3D mineral identification and liberation analysis.  
681 Minerals Engineering, 183, 107592, <https://doi.org/10.1016/j.mineng.2022.107592>.
- 682 Tarquini, S., and Favalli, M. (2010) A microscopic information system (MIS) for petrographic  
683 analysis. Computers & Geosciences, 36, 665–674,  
684 <https://doi.org/10.1016/j.cageo.2009.09.017>.
- 685 Thrun, S., and Mitchell, T.M. (1995) Lifelong robot learning. Robotics and autonomous systems,  
686 15, 25–46, [https://doi.org/10.1016/0921-8890\(95\)00004-y](https://doi.org/10.1016/0921-8890(95)00004-y).
- 687 Torrey, L., and Shavlik, J. (2010) Transfer Learning. In Handbook of Research on Machine  
688 Learning Applications and Trends: Algorithms, Methods, and Techniques pp. 242–264. IGI  
689 Global.
- 690 Trauth, M.H., Gebbers, R., Marwan, N., and Sillmann, E. (2007) MATLAB recipes for earth  
691 sciences Vol. 34. Springer.
- 692 Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., and  
693 Polosukhin, I. (2017) Attention is All you Need. In Advances in Neural Information  
694 Processing Systems Vol. 30. Curran Associates, Inc,  
695 <https://doi.org/10.48550/arXiv.1706.03762>.
- 696 Vyas, A., Katharopoulos, A., and Fleuret, F. (2020) Fast Transformers with Clustered Attention.  
697 In Advances in Neural Information Processing Systems Vol. 33, pp. 21665–21674. Curran  
698 Associates, Inc, <https://doi.org/10.48550/arXiv.2007.04825>.
- 699 Wang, H., Wang, Z., Du, M., Yang, F., Zhang, Z., Ding, S., Mardziel, P., and Hu, X. (2020, April  
700 13) Score-CAM: Score-Weighted Visual Explanations for Convolutional Neural Networks.  
701 arXiv, <https://doi.org/10.48550/arXiv.1910.01279>.



- 702 Wang, Y., Qiu, K.-F., Müller, A., Hou, Z.-L., Zhu, Z.-H., and Yu, H.-C. (2021) Machine learning  
703 prediction of quartz forming-environments. *Journal of Geophysical Research: Solid Earth*,  
704 126, e2021JB021925, <https://doi.org/10.1029/2021JB021925>.
- 705 Wang, Y.D., Shabaninejad, M., Armstrong, R.T., and Mostaghimi, P. (2021) Deep neural networks  
706 for improving physical accuracy of 2D and 3D multi-mineral segmentation of rock micro-  
707 CT images. *Applied Soft Computing*, 104, 107185,  
708 <https://doi.org/10.1016/j.asoc.2021.107185>.
- 709 Wu, M., and Chen, L. (2015) Image recognition based on deep learning. In 2015 Chinese  
710 Automation Congress (CAC) pp. 542–546. Presented at the 2015 Chinese Automation  
711 Congress (CAC).
- 712 Wu, X., Sahoo, D., and Hoi, S.C.H. (2020) Recent advances in deep learning for object detection.  
713 *Neurocomputing*, 396, 39–64, <https://doi.org/10.1016/j.neucom.2020.01.085>.
- 714 Xie, M., Jean, N., Burke, M., Lobell, D., and Ermon, S. (2016) Transfer Learning from Deep  
715 Features for Remote Sensing and Poverty Mapping. *Proceedings of the AAAI Conference on*  
716 *Artificial Intelligence*, 30, <https://doi.org/10.1609/aaai.v30i1.9906>.
- 717 Xing, F., Xie, Y., Su, H., Liu, F., and Yang, L. (2018) Deep Learning in Microscopy Image  
718 Analysis: A Survey. *IEEE Transactions on Neural Networks and Learning Systems*, 29,  
719 4550–4568. Presented at the IEEE Transactions on Neural Networks and Learning Systems,  
720 <https://doi.org/10.1109/tnnls.2017.2766168>.
- 721 Xu, Z., Ma, W., Lin, P., Shi, H., Pan, D., and Liu, T. (2021) Deep learning of rock images for  
722 intelligent lithology identification. *Computers & Geosciences*, 154, 104799,  
723 <https://doi.org/10.1016/j.cageo.2021.104799>.

- 724 Yang, H., and Yang, D. (2023) CSwin-PNet: A CNN-Swin Transformer combined pyramid  
725 network for breast lesion segmentation in ultrasound images. *Expert Systems with*  
726 *Applications*, 213, 119024, <https://doi.org/10.1016/j.eswa.2022.119024>.
- 727 Zhang, L., Qiu, K., Hou, Z., Pirajno, F., Shivute, E., and Cai, Y. (2021) Fluid-rock reactions of the  
728 Triassic Taiyangshan porphyry Cu-Mo deposit (West Qinling, China) constrained by  
729 QEMSCAN and iron isotope. *Ore Geology Reviews*, 132, 104068,  
730 <https://doi.org/10.1016/j.oregeorev.2021.104068>.
- 731 Zhao, Z.-Q., Zheng, P., Xu, S.-T., and Wu, X. (2019) Object Detection With Deep Learning: A  
732 Review. *IEEE Transactions on Neural Networks and Learning Systems*, 30, 3212–3232.  
733 Presented at the *IEEE Transactions on Neural Networks and Learning Systems*,  
734 <https://doi.org/10.1109/tnnls.2018.2876865>.
- 735 Zhong, R., Deng, Y., Li, W., Danyushevsky, L.V., Cracknell, M.J., Belousov, I., Chen, Y., and Li,  
736 L. (2021) Revealing the multi-stage ore-forming history of a mineral deposit using pyrite  
737 geochemistry and machine learning-based data interpretation. *Ore Geology Reviews*, 133,  
738 104079, <https://doi.org/10.1016/j.oregeorev.2021.104079>.
- 739 Zhong, Z., Zheng, L., Kang, G., Li, S., and Yang, Y. (2020) Random Erasing Data Augmentation.  
740 *Proceedings of the AAAI Conference on Artificial Intelligence*, 34, 13001–13008,  
741 <https://doi.org/10.1609/aaai.v34i07.7000>.
- 742 Zhou, B., Khosla, A., Lapedriza, A., Oliva, A., and Torralba, A. (2016) Learning Deep Features  
743 for Discriminative Localization pp. 2921–2929. Presented at the *Proceedings of the IEEE*  
744 *Conference on Computer Vision and Pattern Recognition*,  
745 <https://doi.org/10.1109/cvpr.2016.319>.

746 Zhou, T., Qiu, K., Wang, Y., Yu, H., and Hou, Z. (2022) Apatite Eu/Y-Ce discrimination diagram:  
747 A big data based approach for provenance classification. *Acta Petrol. Sin.*, 38, 291–299,  
748 <https://doi.org/10.18654/1000-0569/2022.01.19>.

749 Zussman, J. and others (1967) Physical methods in determinative mineralogy. *Physical methods*  
750 *in determinative mineralogy*, <https://doi.org/10.1111/j.1365-2818.1979.tb00182.x>.

751

## 752 **Figure Captions**

753 **Figure 1.** The representative images for five minerals which are collected from gold deposits. (a-  
754 d) Gold; (e-h) Pyrite; (i-l) Chalcopyrite; (m-p) Arsenopyrite; (q-t) Stibnite. Apy: arsenopyrite; Au:  
755 gold; Cal: calcite; Ccp: chalcopyrite; Py: pyrite; Qz: quartz; Ser: sericite; Stb: stibnite; Tur:  
756 tourmaline.

757

758 **Figure 2.** Architecture of Swin Transformer. The blocks with different colors represent different  
759 functions. The network has four stages and the last three stages have same structure. The olive  
760 block is patch partition module and patch merging layer. The sage is linear embedding layer. The  
761 light salmon is fully connected layer. And the orange is the classifier. H: height; W: width; C:  
762 color.

763

764 **Figure 3.** Architecture of ResNet50. The blocks with different colors represent different functions.  
765 The network has five stages and the last four stages have same structure. Stage 1: the dark salmon  
766 block is the convolutional layer. the salmon is the normalization. the light one is the activation  
767 function. and the sage one is the max pooling layer; Stage 2 to stage 5: olive one is the  
768 convolutional block with one convolutional layer; the salmon block is the identity block with two,

769 three, five and two convolutional layers, respectively. And the last light salmon is the pooling layer,  
770 and the orange is the classifier.

771

772 **Figure 4.** Architecture of MobileNetv2. It contains two units: stride=1 and stride=2. Conv 1x1 is  
773 the 1x1 convolutional kernel. ReLU is nonlinear activation function. Dwise 3x3 is depth-wise  
774 convolution with 3x3 convolutional kernel.

775

776 **Figure 5.** Workflow of the proposed automatic classification. Step 1: dataset compiling. Crop raw  
777 images using OpenCV. Select the processed images that contain only one mineral phase; Step 2:  
778 data splitting and augmentation. The dataset was divided into training set, validation set and test  
779 set (3:1:1). The data augmentation methods include random erasing, flipping, brightness adjust,  
780 random zoom, random contrast and random saturation; Step 3: model training and evaluating. Swin  
781 Transformer, ResNet50 and MobileNetv2 algorithms were used to train the classification models.  
782 The model evaluation metrics include accuracy, precision, recall, and F1-score; Step 4: model  
783 predicting. Put the images to the trained model to predict the five metal mineral classes.

784

785 **Figure 6.** Changes of (a) validation accuracy and (b) training loss of three algorithms using the  
786 method of transfer learning. The lines reflect the changes of different algorithms' performance  
787 within 20 epochs (green: Swin Transformer algorithm; red: ResNet50 algorithm; blue:  
788 MobileNetv2 algorithm).

789

790 **Figure 7.** Confusion matrix of the test set used to evaluate the three algorithms. (a) Swin  
791 Transformer; (b) ResNet50; (c) MobileNetv2. Indicated values are the number of images. The

792 horizontal axis represents the predicted label, while the vertical axis denotes the true label. The  
793 horizontal axis is the predicted label, while the vertical axis is the true label. Apy: arsenopyrite;  
794 Ccp: chalcopyrite; Au: gold; Py: pyrite; Stb: stibnite.

795

796 **Figure 8.** Presentation of erroneous classification results from Swin Transformer metal mineral  
797 classifier. The model misclassified (a) pyrite and (b) chalcopyrite, (c) arsenopyrite and (d) stibnite,  
798 as well as (e) pyrite and (f) arsenopyrite. Apy: arsenopyrite; Ccp: chalcopyrite; Py: pyrite; Stb:  
799 stibnite.

800

801 **Figure 9.** Changes in (a) validation accuracy and (b) training loss of Swin Transformer with  
802 transfer learning and without transfer learning respectively. The lines reflect the changes of  
803 different algorithms' performance within 20 epochs (dark green: Swin Transformer with transfer  
804 learning; light green: Swin Transformer without transfer learning).

805

806 **Figure 10.** CAMs of three models with five-classes metal minerals image classification. The  
807 redder the mapping, the higher the response of the corresponding area of the original image to the  
808 model's classification output. The numbers on the mappings represent the output probability of  
809 the model for unknown minerals. Ccp: chalcopyrite; Py: pyrite; Au: gold; Apy: arsenopyrite; Stb:  
810 stibnite.

811

812

813

814

815

816

817

818

819

820

821 **Table**822 **Table 1** Characteristics of Metal Minerals Under the Microscope

	gold	pyrite	chalcopyrite	arsenopyrite	stibnite
<b>Chemical composition</b>	Au	FeS <sub>2</sub>	CuFeS <sub>2</sub>	FeAsS	Sb <sub>2</sub> S <sub>3</sub>
<b>Reflectivity</b>	Gold 480: 33.97; 546: 70.67; 589: 80.09; 656: 85.88	White: 54.5; 470: 46; 546: 53; 589: 54; 650: 54	White: 44–46.1; 470: 34; 546: 47; 589: 48; 650: 49	White: 51.7–55.7; 470: 51–55; 546: 52–54; 589: 53–54; 650: 53	White: 30.2–40; 470: 31–53; 546: 31–48; 589: 30–45; 650: 30–42
<b>Reflection color</b>	Golden yellow, bright yellow	Light yellow	Copper yellow	Bright white with cream or red color, weak polychromatic	White to light off-white
<b>Homogeneity and heterogeneity</b>	Homogenous	Homogenous	Weak heterogeneity	Strong heterogeneity	Strong heterogeneity
<b>Morphological characteristics</b>	Polymeric crystals between octahedral, hexahedral, tetrahedral, triangular, and rhomboid dodecahedron; Irregularly granular	Euhedral crystals in the form of cubes, pentagonal dodecahedron and octahedron	Irregular granular crystals	Diamond-shaped, elongated columnar, spear-headed and other euhedral crystals	Columnar long and granular crystals
<b>Mineral combination</b>	Arsenopyrite, pyrite, chalcopyrite, pyrrhotite, galena, sphalerite, stibnite, calcite, tellurite and quartz	Iron, copper, lead, zinc, silver sulfide, gold, rutile, graphite, etc.	Associated with sulfides	Pyrite, loellingite, tetrahedrite, magnetite, galena, sphalerite, stibnite, etc.	Berthierite, pyrite, arsenopyrite, sphalerite, tetrahedrite, scheelite, gold, realgar, orpiment, etc.
<b>References</b>	(Shang and Lin 1990; Piller 2012; Ramdohr 2013)	(Cameron 1961; Chen et al. 1979; Ramdohr 2013)	(Piller 1966; Zussman and others 1967; Santosh et al. 2007)	(Picot and Johan 1977; Lu and Peng 2010)	(Craig et al. 1981; Criddle and Stanley 2012)

823

824

825

826

827

828 **Table 2** Summary of the Training, Validation and Test Sets of Image Dataset

	Arsenopyrite	Chalcopyrite	Gold	Pyrite	Stibnite	Total
<b>Training</b>	548	505	490	574	596	2713
<b>Validation</b>	182	170	164	191	198	905
<b>Test</b>	183	169	163	192	199	906
<b>Total</b>	913	844	817	957	993	4524

829

830

831

832 **Table 3** Mineral Classification Performance on the Test Set

Method	Metric	Arsenopyrite	Chalcopyrite	Gold	Pyrite	Stibnite	Metric Value
	Acc						<b>0.92</b>
Swin-Transformer	Pre	0.85	0.83	0.96	1.00	0.99	<b>0.93</b>
	Rec	0.99	0.97	1.00	0.81	0.84	<b>0.92</b>
	F1	0.91	0.89	0.98	0.90	0.91	<b>0.92</b>
	Acc						<b>0.91</b>
ResNet50	Pre	0.97	0.78	0.88	0.94	0.97	<b>0.91</b>
	Rec	0.98	0.89	0.99	0.70	0.98	<b>0.91</b>
	F1	0.97	0.84	0.93	0.80	0.98	<b>0.90</b>
	Acc						<b>0.81</b>
MobileNetv2	Pre	0.64	0.77	0.91	0.98	1.00	<b>0.86</b>
	Rec	1.00	0.90	0.99	0.61	0.62	<b>0.82</b>
	F1	0.78	0.83	0.95	0.75	0.76	<b>0.81</b>

833 **Note:** Acc: abbreviation for model evaluation indicator accuracy; Pre: abbreviation for model evaluation indicator precision; Rec:

834 abbreviation for model evaluation indicator recall; F1: abbreviation for model evaluation indicator F1-score.

Figure 1

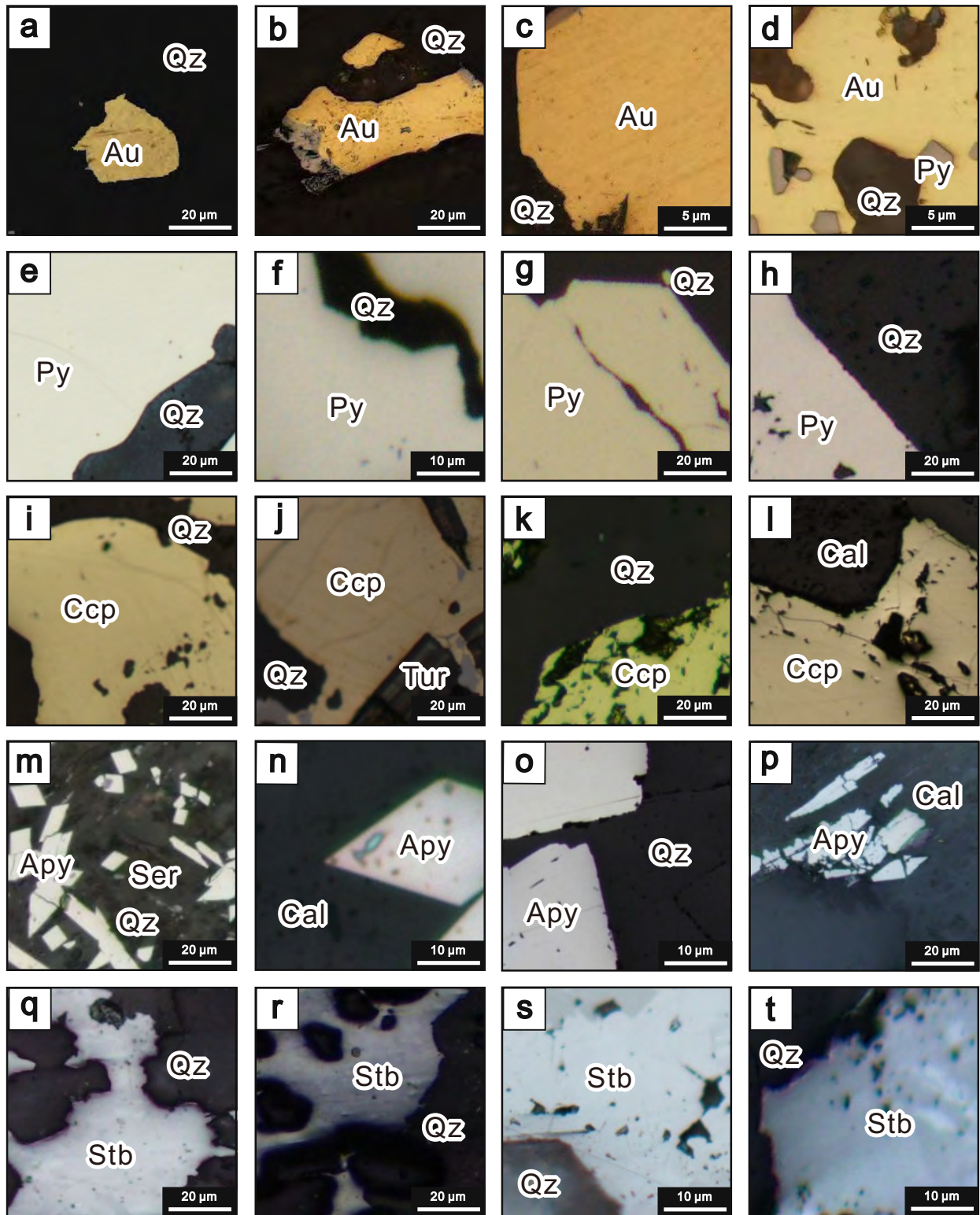




Figure 2

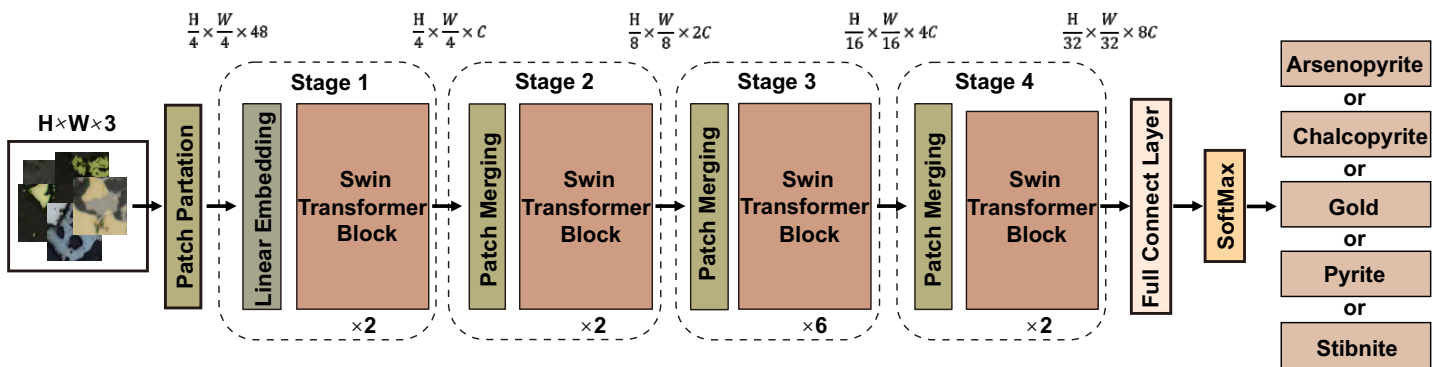


Figure 3

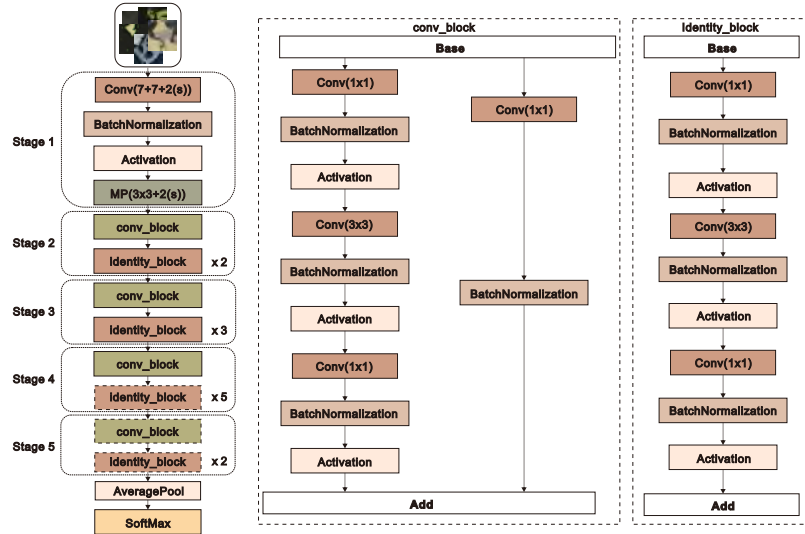


Figure 4

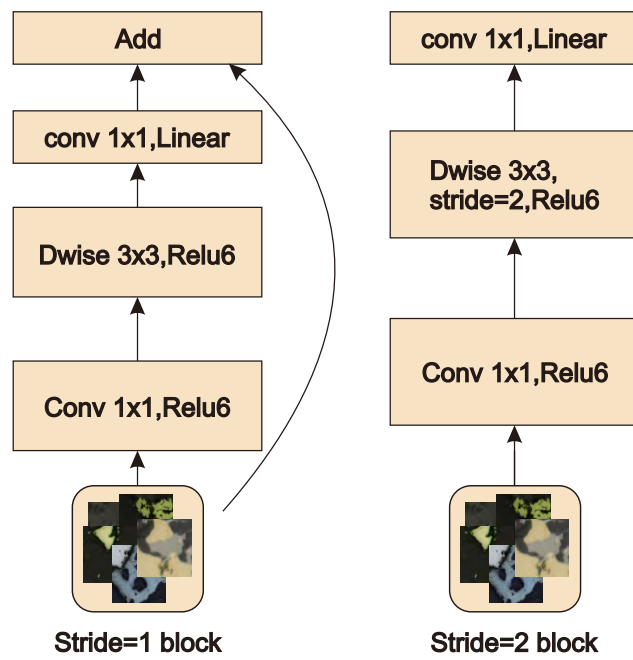


Figure 5

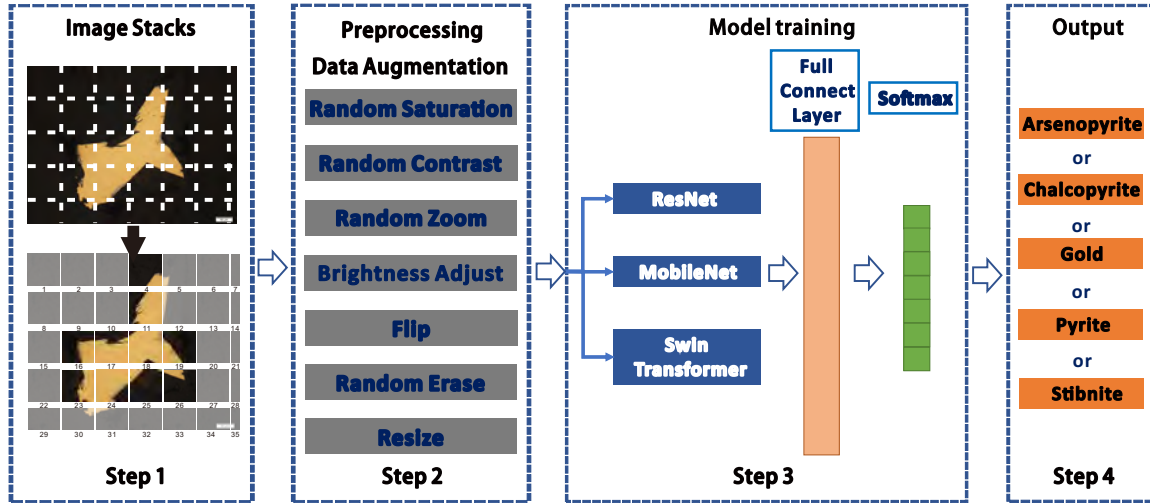


Figure 6

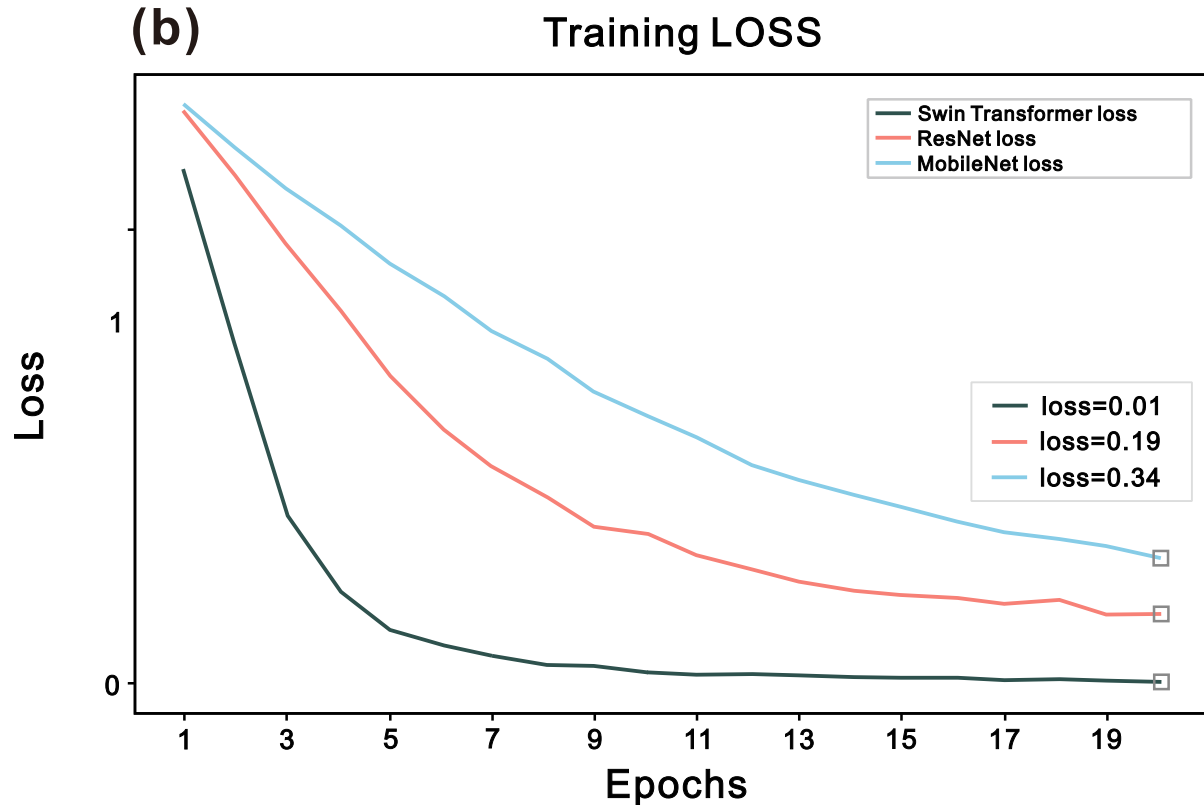
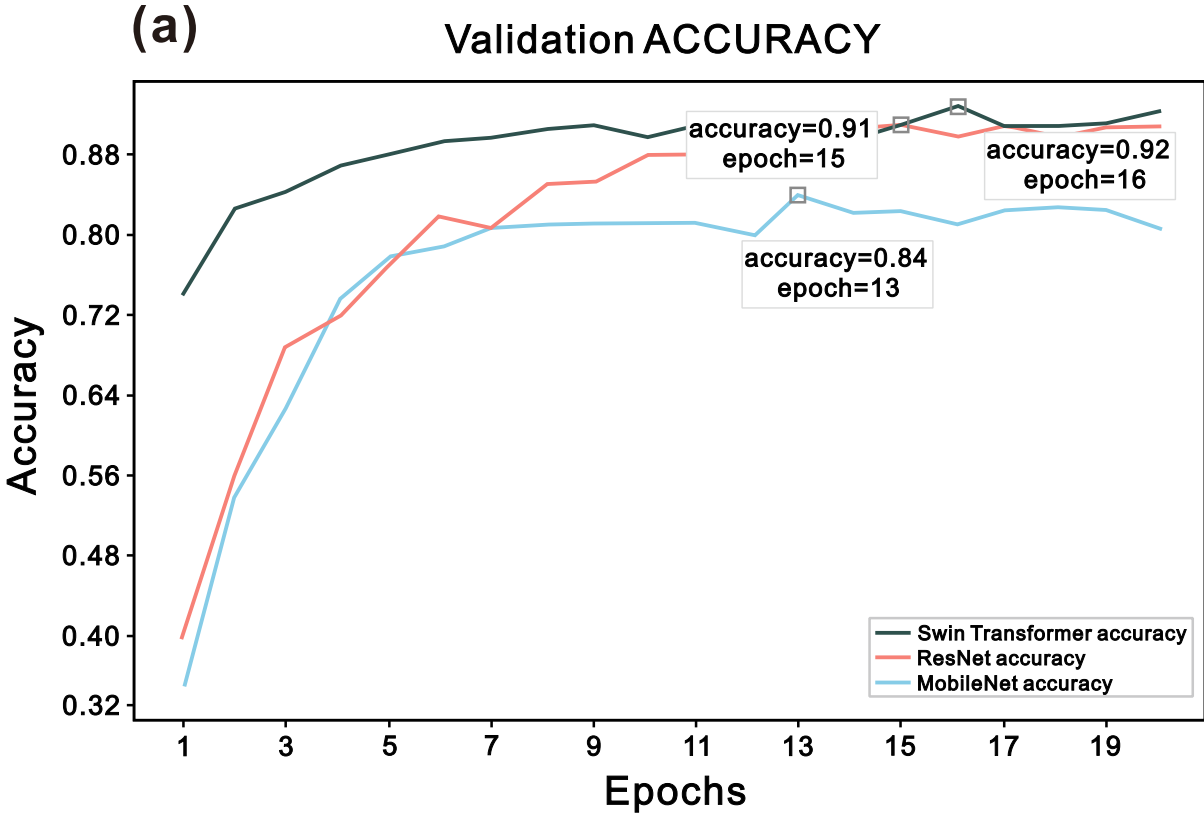


Figure 7

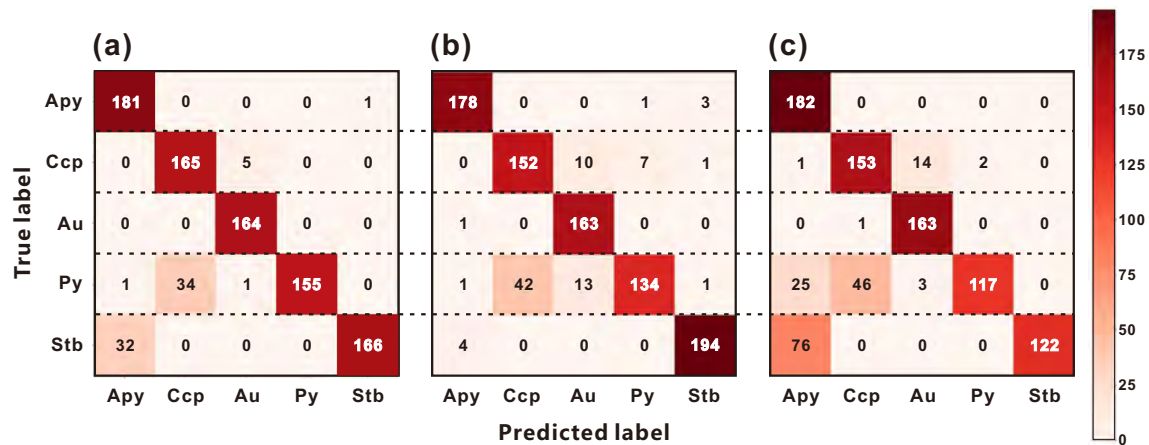


Figure 8

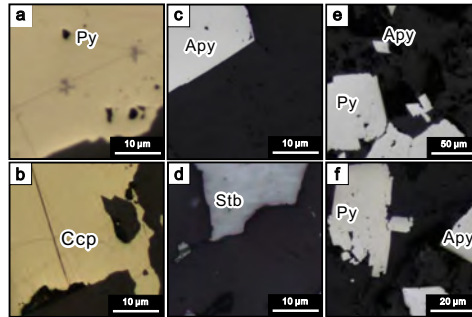


Figure 9

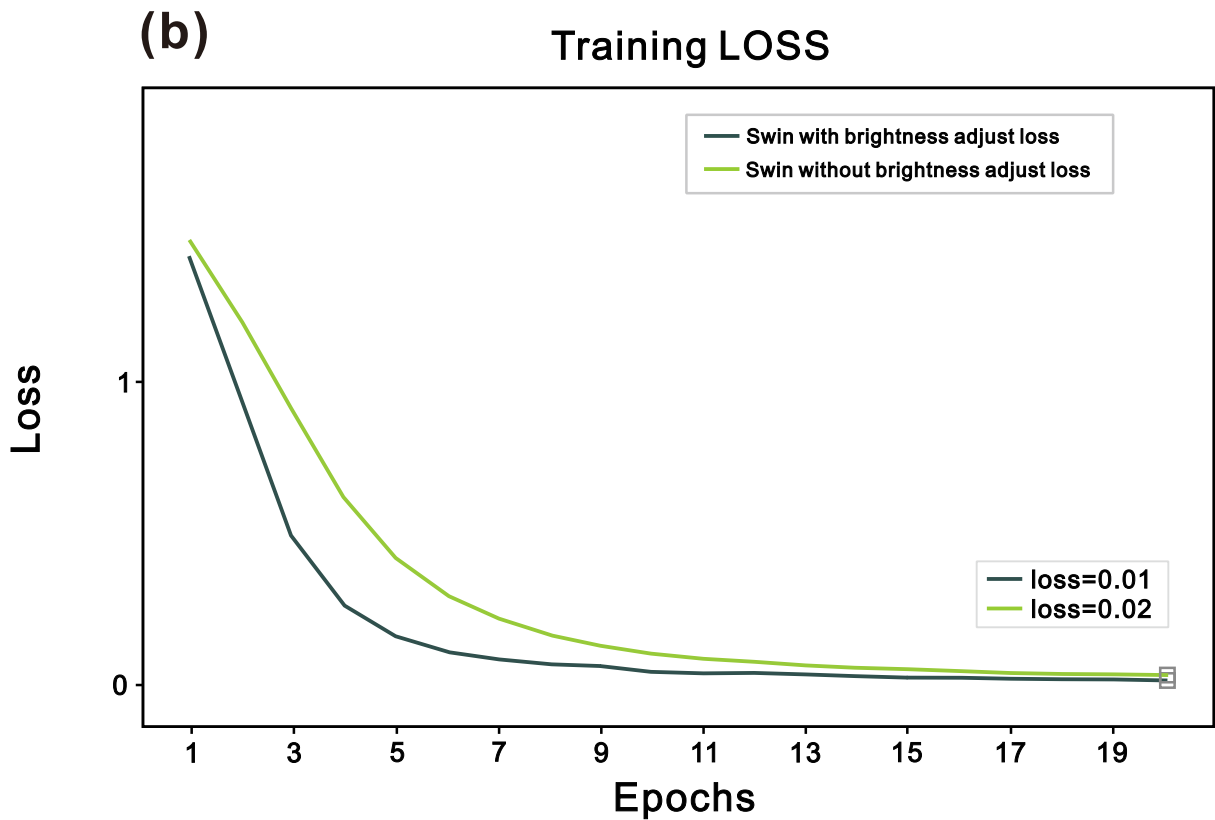
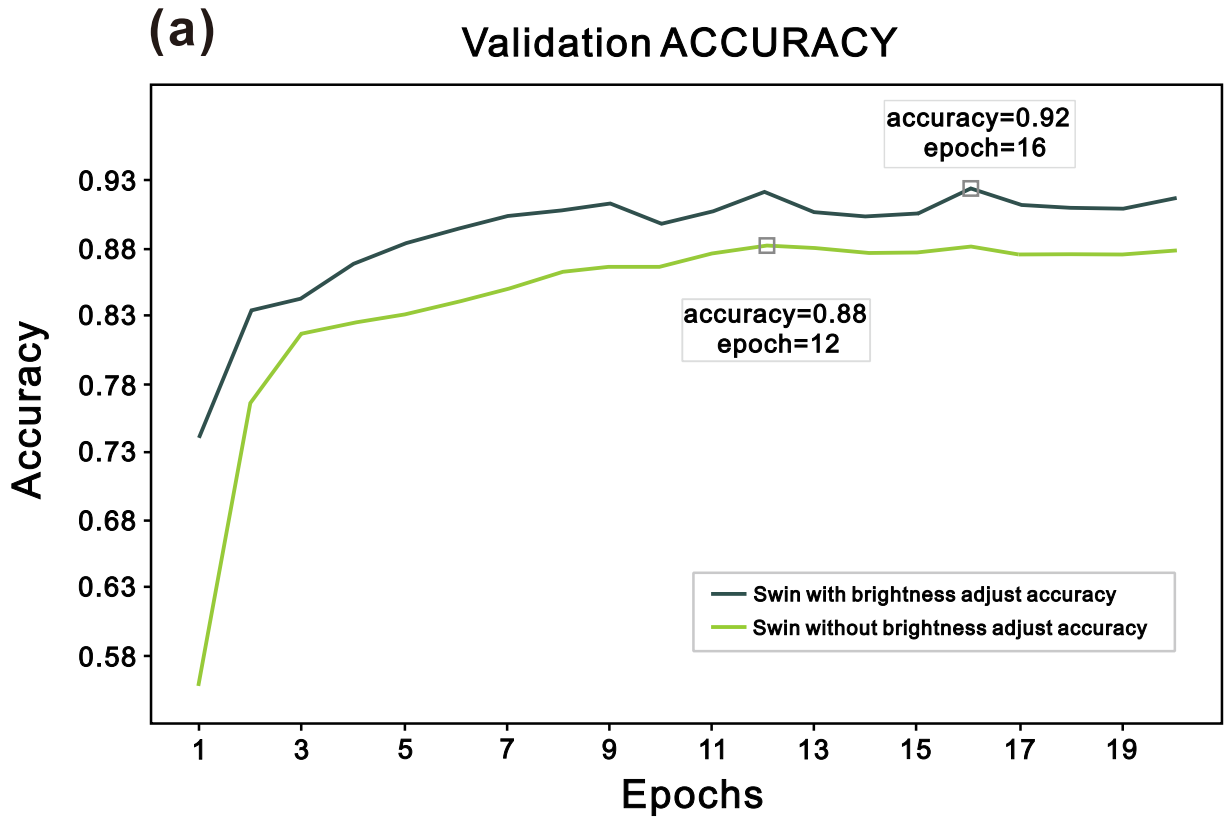




Figure 10

