# CLASSIFICATION OF THE CERVICAL VERTEBRAE MATURATION (CVM) STAGES USING THE TRIPOD NETWORK

*Salih Atici*[†‡]    *Hongyi Pan*[†‡]    *Mohammed H. Elnagar**    *Veerasathpurush Allareddy**

*Omar Suhaym**[⋆]    *Rashid Ansari*[†]    *Ahmet Enis Cetin*[†]

[†]Department of Electrical and Computer Engineering, University of Illinois Chicago, Chicago, IL
*Department of Orthodontics, College of Dentistry, University of Illinois Chicago, Chicago, IL
[⋆]Department of Oral and Maxillofacial Surgery, King Saud bin Abdulaziz
University for Health Sciences, Riyadh, Saudi Arabia

## ABSTRACT

We present a novel deep learning method for fully automated detection and classification of the Cervical Vertebrae Maturation (CVM) stages. The deep convolutional neural network consists of three parallel networks (TriPodNet) independently trained with different initialization parameters. They also have a built-in set of novel directional filters that highlight the Cervical Verte edges in X-ray images. Outputs of the three parallel networks are combined using a fully connected layer. 1018 cephalometric radiographs were labeled, divided by gender, and classified according to the CVM stages. Resulting images, using different training techniques and patches, were used to train TripodNet together with a set of tunable directional edge enhancers. Data augmentation is implemented to avoid overfitting. TripodNet achieves the state-of-the-art accuracy of 81.18% in female patients and 75.32% in male patients. The proposed TripodNet achieves a higher accuracy in our dataset than the Swin Transformers and the previous network models that we investigated for CVM stage estimation.

*Index Terms*— Deep learning, cervical vertebrae maturation, tripod network, vision transformers.

## 1. INTRODUCTION

The success of orthodontic/orthopedic treatment depends on optimal treatment timing, especially in addressing craniofacial skeletal imbalances. The optimal treatment timing relies on identifying craniofacial skeletal maturity stages. Bone age assessment using radiographic analyses was reported to be more accurate than chronological age in determining skeletal maturation, growth rate, the peak period of growth, and the remaining growth potential [1, 2, 3]. Cervical vertebra maturation (CVM) staging in lateral cephalometric radiographs is a method to determine skeletal maturation. The validity and reliability of the CVM staging have been supported by multiple studies [4, 5]. Cervical vertebrae are the first seven bones of the spinal column. Vertebral growth involves changes in the size of vertebral bodies and the shape of the upper and lower borders of C2, C3, C4 vertebrae. These changes have been described into six stages, correlating with morphological modifications of the vertebral shapes. The major limitation of the CVM method is that it

is not user-friendly and needs experienced practitioners; researchers reported poor reproducibility among nonexpert examiners [6].

The use of machine learning (ML) techniques in the field of medical imaging is rapidly evolving, and a fully automated diagnostic approach has gained attention with its promise of reducing human error as well as the time and effort needed for the task [7]. The application of Deep Learning (DL) to study human growth and development from radiographs is a promising idea that needs to be explored. The present study aims to apply a custom-designed DL method to develop a fully automated machine-learning system to detect and classify the CVM stages. There have been recent studies to use pre-trained networks to create fully automated detection and classification of the CVM stages where each utilizes a different dataset [8, 9]. In this study, we propose a custom-designed network model to develop a fully automated system to detect and classify the CVM stages. Our DL network has a tripod-like structure consisting of three parallel networks which are independently trained with different initialization parameters [10]. They also have a built-in set of novel directional filters that highlight the edges of the cervical vertebrae in X-ray images.

The TripodNet has some of the features of the transformer networks: (i) Input images are divided into patches, (ii) input patches are augmented as in transformers; and (iii) the network has a multi-headed structure. ResNet-20 is used as the backbone of TripodNet. Output feature maps of the three parallel ResNets are combined using a fully connected layer to produce the final decision as shown in Fig. 1. Moreover, the age information of the patients is also fed to the network to increase the accuracy of classification. The dataset is divided by gender as male patients can have a different growth rate than female patients. The proposed model achieves state-of-the-art performance with 81.18% in female patients and 75.32% in male patients on the dataset collection, which is superior to the previous results that used DL in a traditional way including the straightforward implementation of the vision transformer [11].
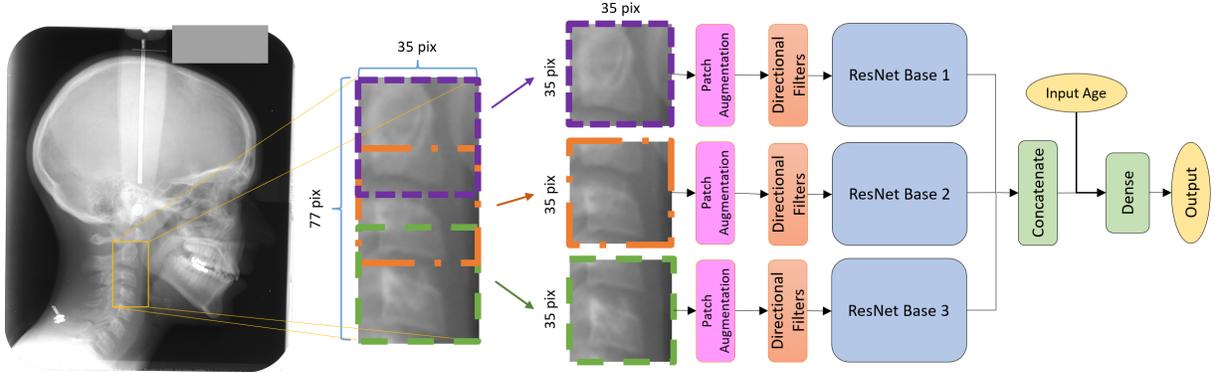
## 2. METHODOLOGY

### 2.1. Overview of the CVM Dataset

The dataset used in developing our algorithm consists of digitized images of scanned lateral cephalometric films for subjects aged between 4 and 29 obtained from the American Association of Orthodontists Foundation (AAOF) Craniofacial Growth Legacy Collections, an open data source [12]. The images were studied and labeled by the third author (MHE) who is an expert orthodontist scien-
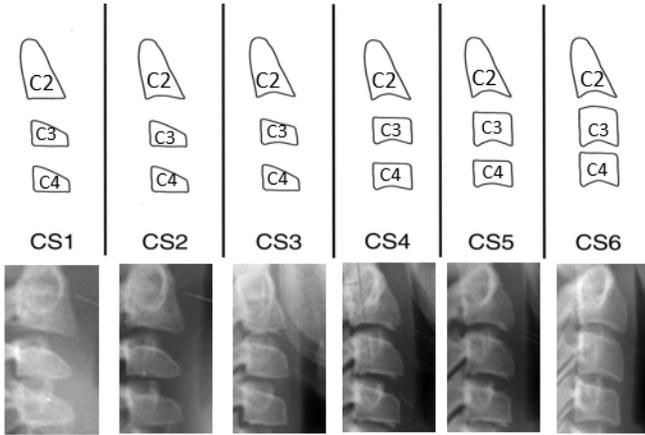
**Fig. 1**: The Model Diagram of the TriPodNet. Overlapping vertebrae images are fed to each pod of the network.

tist with more than ten years of experience in classifying CVM. Cervical maturation stages were classified into six stages (CS1- CS6). Six stages are shown in Fig. 2 [2]. It is visible from Fig. 2 that the main difference among the classes arises in the size and shape of C2, C3, and C4 vertebrae. The change may also happen in one vertebrae only, which may cause confusion for a traditional CNN model which uses the entire image as its input. For example, two classes CS1 and CS2 are only separated by the shape of the lower bound of C2 vertebra whereas C3 and C4 vertebrae have the same shape.



**Fig. 2**: Six CVM Stages and the corresponding X-ray images. The top drawing is adapted from reference [2].

The main dataset was classified into six stages of CVM (CS1-CS6) and we will examine the model performances on the six-stage classification problem. The dataset consists of 1012 images classified by the third author (MHE) as the principal evaluator. Out of 1012 images, 478 images belong to female patients, and 534 images belong to male patients. In our entire dataset, the number of lateral cephalograms belonging to cervical stages CS1, CS2, CS3, CS4, CS5, and CS6 are 153, 182, 174, 159, 167, and 177, respectively. Data augmentation methods such as random translation, AutoContrast [13], AugMix [14], and RandAugment [15] are implemented with various combinations to avoid the overfitting of the DL networks to the training dataset.

As it is stated earlier, the changes in C2, C3, and C4 vertebrae create 6 classes; therefore, their shape and size play a significant role in the determination of CVM stages. TriPodNet requires three inputs

to create an output. In this study, for the network to distinguish the differences among the stages, we used three patches per image where each patch containing one of the C2, C3, or C4 vertebrae. Moreover, every patch is augmented at the beginning of the TriPodNet to increase the performance of each model. Since the size of the dataset is much smaller than the benchmark dataset used to train ResNet-20, data augmentation and patch augmentation methods are necessary to train the TriPodNet. Note that the patch augmentation is different from the data augmentation method we used to avoid overfitting. Rotation and grayscale jittering are implemented to create a different instance of the same input each time an input is about to be fed into the network.

## 2.2. Structure of the Multi-Pod Network (MultiPodNet)

A typical MultiPodNet utilizes two or more parallel Convolutional Neural Networks (CNNs) performing the same computations and may process the input image as sequential image patches in parallel as in transformer networks. The original input image and its augmented versions are fed into convolutional networks forming the MultiPod network and the output feature maps of parallel convolutional networks are concatenated before the fully connected dense layer. In this study, we use TriPodNet because it generated the best results in our dataset and it is the best performer among the other models compared to our arXiv manuscript [10]. Similar to [10], we also use ResNets as the baseline network in this article. The TriPodNet model structure is shown in Fig 1.

We first segment a given head and neck X-ray image and identify the spine (cervical vertebrae) region using the so-called Aggregate Channel Features (ACF) object detector [16] as shown in Fig. 1. This avoids the process of manually cropping the spine region in each image in the database. As a result, the skull, jaw, and irrelevant background regions are removed before the images are applied to the deep learning algorithm. The ACF object detector automatically extracts the Region of Interest (RoI) in the images thereby reducing the search space of the deep learning structure. Because all the segmented images have variable sizes, they are resized to a common size of $77 \times 35$.

After this step, we use image patches to create input image sequences for the TriPodNet model. Shapes of C2, C3 and C4 vertebrae determine the stage of the CVM. Therefore, we crop the segmented vertebrae image into $35 \times 35$ patches that contain only one vertebra as shown in Fig. 1. Three patches are derived from each image before the patch augmentation. The location of each vertebra in the image is used to create the patches since C2, C3 and C4 are

always in similar positions in any image. We used overlapped windows to create patches with the size of $35 \times 35$. Before using the patches in model training, we also benefit from patch augmentation where we rotate the images randomly by 5 degrees. The rotation is necessary to help the model generalize as the alignment of each vertebra depends on the posture of the patient. Moreover, we apply grayscale jittering on the patches. With the grayscale jittering, patches become brighter or darker, randomly. Patch augmentation is visualized in Fig. 3. Similar to other image classification methods, augmentation is only applied to the training images. We do not apply patch augmentation on testing images.
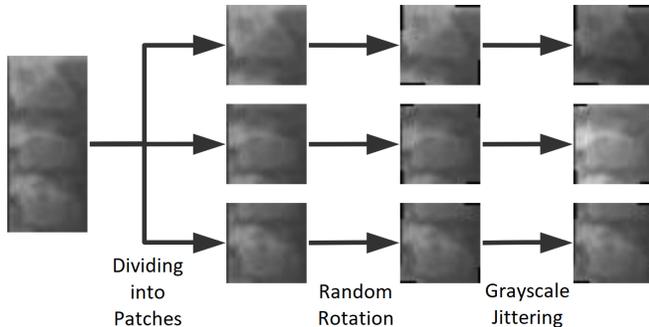


**Fig. 3**: Patch augmentation on the training images.

Next, instead of feeding the image patches directly to the ResNets, the edges of the vertebral body in the patches are emphasized using eight directional filters described in [17, 18]. The outputs of directional filters are then fed to the ResNets. The same directional filters in [11] are used before ResNets to highlight the edges of CVM images. The motivation behind using the directional filters is to start the deep learning model using our prior domain knowledge. Since the angles of CVM bones are the distinguishing factor in estimating the CVM stages highlighting the edges of the bones in multiple directions will give an advantage to the deep network model.

As pointed out above each pod of the TriPodNet is a ResNet-20. The structure of ResNet-20 with the directional filters and patch augmentation is summarized in Table 1. The ResNets are implemented in parallel. The output feature maps of ResNets are concatenated together with the age of the subject. We also use the chronological age information as the input to the fully connected layer. Obviously, the chronological age of a subject is correlated with the maturity of a patient. The age information is repeated six times in the vector and Gaussian noise with zero mean and 0.01 variance is added to secure its impact before the output layer.

We selected the TriPodNet with random translation and AutoContrast augmentation methods [13] as the best network model in our dataset based on our experiments. In the next section, we present our experimental results.

## 3. EXPERIMENTAL RESULTS

We studied different networks, and different data augmentation methods to determine the best possible network structure. As introduced in [10], MultiPodNet can be constructed from two or more pod networks processing the input in parallel. Outputs of pod networks can be combined by adding the feature maps or by concatenation that produced a higher accuracy than adding the feature maps of individual networks. We studied a single channel ResNet,

| Layer | Output Shape | Implementation Details |
|---|---|---|
| PatchAug | $35 \times 35 \times 8$ | - |
| DirFilts | $35 \times 35 \times 8$ | $7 \times 7, 8$ |
| Conv1 | $35 \times 35 \times 16$ | $3 \times 3, 16$ |
| Conv2_x | $35 \times 35 \times 16$ | $\begin{bmatrix} 3 \times 3, 16 \\ 3 \times 3, 16 \end{bmatrix} \times 3$ |
| Conv3_x | $17 \times 17 \times 32$ | $\begin{bmatrix} 3 \times 3, 32 \\ 3 \times 3, 32 \end{bmatrix} \times 3$ |
| Conv4_x | $8 \times 8 \times 64$ | $\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 3$ |
| AAP | 64 | Adaptive Average Pooling |

**Table 1**: Structure of each pod (ResNet with the directional filters) of the TriPodNet. PatchAug stands for patch augmentation. DirFilts stands for directional filters [11].

DuPodNet which consists of two pod networks, TriPodNet, and QuadPodNet which consists of four parallel pod networks. The TriPodNet with input overlapping input image patching strategy as shown in Fig. 1 produced the highest accuracy in our dataset. We compare the TripodNet constructed from ResNet-20s [19] with the Swin Transformer [20], Xception [21], MobileNet-V1 [22], MobileNet-V2 [23], and the custom designed CNN with the directional filters that achieved the previous best result in [11]. In addition to random translation and AutoContrast data augmentation methods we studied "randAugment", "AugMix" methods and without any augmentation.

To train the networks, an SGD optimizer with a weight decay of 0.0001 and momentum of 0.9 is used. These models are trained with a batch size of 32, an initial learning rate of 0.1 for 100 epochs, and the learning rate is reduced to 1/10 at epochs 25, 50, and 75. The experiments are implemented using PyTorch in Python 3.

The accuracy results of the TripodNet with and without directional filters and with different augmentation methods are summarized in Table 2. Directional filters improve TripodNet's accuracy by 4.88% (from 76.29% to 81.17%) in the dataset containing female subjects. Similarly, they improve the accuracy by 4.17% (from 71.15% to 75.32%) in the dataset containing male subject images. As pointed out above we augment both the entire image and the patches. We trained the network using random translations with AutoContrast as the data augmentation method as shown in Table 3. Furthermore, we augment the input image patches to make the system robust to changes in posture and exposure as shown in Table 2 during training.

| Directional Filters | Augmentation | | Accuracy | |
|:---:|:---:|:---:|:---:|:---:|
| | Data | Patch | Female | Male |
| × | × | × | 67.05% | 65.38% |
| × | ✓ | × | 75.11% | 70.19% |
| ✓ | ✓ | × | 78.64% | 70.19% |
| × | ✓ | ✓ | 76.29% | 71.15% |
| ✓ | ✓ | ✓ | **81.17%** | **75.32%** |

**Table 2**: Accuracy results of the TriPodNet with and without directional filters, augmentation of entire input images, and augmentation of image patches.
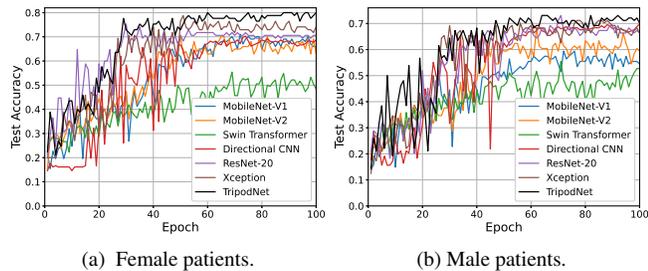
| Data Augmentation Method | Female | Male |
|---|---|---|
| No augmentation | 68.23% | 66.34% |
| Random augmentation | 74.11% | 70.18% |
| AugMix | 77.64% | 69.23% |
| **Random translation, AutoContrast** | **81.17%** | **75.32%** |

**Table 3**: Accuracy results of various data augmentation methods. Directional filters and patch augmentation are applied. AugMix, AutoContrast are PyTorch commands [13].

The accuracy of different MultiPod networks is presented in Table 4. TriPodNet achieves the highest accuracy. QuadPodNet has more parameters compared to TriPodNet but it produced a lower accuracy. This may be due to the small dataset size and training issues. StackNet in Table 4 has also three parallel pod networks. In the StackNet, we stack the input patches and feed the augmented stack patches to the pod networks at the same time.

| Model | Parameters | Female | Male |
|---|---|---|---|
| ResNet | 0.27M | 75.29% | 73.07% |
| DuPodNet | 0.54M | 75.29% | 68.27% |
| StackNet | 0.82M | 80.03% | 70.19% |
| **TriPodNet** | 0.81M | **81.17%** | **75.32%** |
| QuadPodNet | 1.08M | 76.84% | 71.43% |

**Table 4**: Accuracy results of MultiPod networks. DuPodNet uses two ResNet pods, TriPodNet uses three, and QuadPodNet uses four. StackNet feeds all three patches stacked to three pods.



(a) Female patients.  (b) Male patients.

**Fig. 4**: Test accuracy versus the number of epochs of different networks. The TripodNet reaches the highest accuracy on both genders.
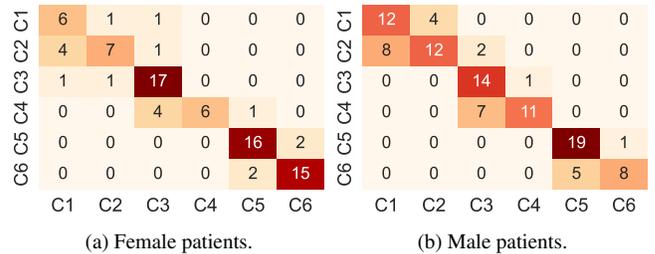
Comparison of the TripodNet with other state-of-art networks is presented in Table 5 and Fig. 4. Fig. 4 shows the test accuracy versus the number of epochs of different networks. In Table 5, we compare the TripodNet with ResNet-20 [19], Swin Transformer [20], Xception [21], MobileNet-V1 [22], MobileNet-V2 [23], and the custom-designed CNN with the directional filters [11]. We adjust the image sizes to pre-trained networks: we remove some of the downsampling layers in the MobileNet-V1, MobileNet-V2, and Xception. We implement the Swin Transformer using two stages, where each stage contains a patch merging layer with 4 Swin Transformer blocks. Moreover, it is well-known that after pre-training on a large dataset, a transformer can get a much higher accuracy (this is known as transfer learning). Therefore, to compare with the Swin Transformer, we

also use the ImageNet-1k-pre-trained Swin-s and Swin-t transformers. In these cases, the input tensors are interpolated to $224 \times 224$ to adjust the CVM image sizes to transformer input sizes. Table 5 and Fig. 4 show that our proposed TripodNet produces the highest accuracy in our data set. The TripodNet gets better results than the pre-trained Swin-s and Swin-t transformers. This may be due to the "tiny" size of the CVM data set. It is probably too small for the transformer networks. Another reason is that transfer learning may not work properly because the CVM images are quite different from the ImageNet images.

| Model | Parameters | Female | Male |
|---|---|---|---|
| MobileNet-V1 | 3.21M | 73.20% | 60.40% |
| MobileNet-V2 | 2.29M | 72.16% | 68.32% |
| Swin Transformer | 4.12M | 63.58% | 62.50% |
| Pre-trained Swin-s | 48.8M | 75.26% | 73.27% |
| Pre-trained Swin-t | 28.5M | 72.16% | 74.26% |
| Directional CNN [11] | 0.71M | 70.33% | 71.88% |
| ResNet-20 | 0.27M | 75.29% | 73.07% |
| Xception | 33.0M | 78.82% | 71.15% |
| **TripodNet** | 0.81M | **81.17%** | **75.32%** |

**Table 5**: Accuracy results of various networks. Swin-s and Swin-t were pre-trained on ImageNet-1K.

Fig. 5 shows the confusion matrices of the TriPodNet in the male and female patient datasets.



(a) Female patients.  (b) Male patients.

**Fig. 5**: The confusion matrices of the TriPodNet.

## 4. CONCLUSION

In this paper, we present a new method for CVM classification. We introduce a novel neural network which is a combination of three parallel networks for this purpose. Similar to the transformer networks, TriPodNet performs its computations in parallel using traditional CNNs and the output feature maps of parallel CNNs are combined using a fully connected layer to produce the final result. We also compared the results of two, three, four, or more parallel networks. The TriPodNet with three parallel ResNet-20s produced the best accuracy result.

The transformer networks did not produce as good results as the TriPodNet. This is probably because our tiny data set is too small for the transformer networks.

## 5. REFERENCES

[1] Sunjay Suri, Chandrakala Prasad, Bryan Tompson, and Wendy Lou, "Longitudinal comparison of skeletal age determined by the greulich and pyle method and chronologic age in normally growing children, and clinical interpretations for orthodontics," *American Journal of Orthodontics and Dentofacial Orthopedics*, vol. 143, no. 1, pp. 50–60, 2013.

[2] Lorenzo Franchi, Tiziano Baccetti, Laura De Toffol, Antonella Polimeni, and Paola Cozza, "Phases of the dentition for the assessment of skeletal maturity: a diagnostic performance study," *American Journal of Orthodontics and Dentofacial Orthopedics*, vol. 133, no. 3, pp. 395–400, 2008.

[3] Zachary J Mellion, Rolf G Behrents, and Lysle E Johnston Jr, "The pattern of facial skeletal growth and its relationship to various common indexes of maturation," *American Journal of Orthodontics and Dentofacial Orthopedics*, vol. 143, no. 6, pp. 845–854, 2013.

[4] Leonard S Fishman, "Radiographic evaluation of skeletal maturation: a clinically oriented method based on hand-wrist films," *The Angle Orthodontist*, vol. 52, no. 2, pp. 88–112, 1982.

[5] Brent Hassel and Allan G Farman, "Skeletal maturation evaluation using cervical vertebrae," *American Journal of Orthodontics and Dentofacial Orthopedics*, vol. 107, no. 1, pp. 58–66, 1995.

[6] Trenton S Nestman, Steven D Marshall, Fang Qian, Nathan Holton, Robert G Franciscus, and Thomas E Southard, "Cervical vertebrae maturation method morphologic criteria: poor reproducibility," *American Journal of Orthodontics and Dentofacial Orthopedics*, vol. 140, no. 2, pp. 182–188, 2011.

[7] June-Goo Lee, Sanghoon Jun, Young-Won Cho, Hyunna Lee, Guk Bae Kim, Joon Beom Seo, and Namkug Kim, "Deep learning in medical imaging: general overview," *Korean journal of radiology*, vol. 18, no. 4, pp. 570–584, 2017.

[8] Hyejun Seo, JaeJoon Hwang, Taesung Jeong, and Jonghyun Shin, "Comparison of deep learning models for cervical vertebral maturation stage classification on lateral cephalometric radiographs," *Journal of Clinical Medicine*, vol. 10, no. 16, pp. 3591, 2021.

[9] Hossein Mohammad-Rahimi, Saeed Motamadian, Mohadeseh Nadimi, Sahel Hassanzadeh-Samani, Mohammad Minabi, Erfan Mahmoudinia, Victor Lee, and Mohammad Hossein Rohban, "Deep learning for the classification of cervical maturation degree and pubertal growth spurts: A pilot study," *Korean journal of orthodontics*, vol. 52, pp. 112–122, 03 2022.

[10] Hongyi Pan, Salih Atici, and Ahmet Enis Cetin, "Multipod convolutional network," *arXiv preprint arXiv:2210.00689*, 2022.

[11] Salih Furkan Atici, Rashid Ansari, Veerasathpurush Allareddy, Omar Suhaym, Ahmet Enis Cetin, and Mohammed H Elnagar, "Fully automated determination of the cervical vertebrae maturation stages using deep learning with directional filters," *Plos one*, vol. 17, no. 7, pp. e0269198, 2022.

[12] "Aaof craniofacial growth legacy collection," `https://www.aaoflegacycollection.org/aaof_home.html` Accessed: 2021-10-03.

[13] "Pytorch autocontrast," `https://pytorch.org/vision/stable/generated/torchvision.transforms.functional.autocontrast.html` Accessed: 2022-10-03.

[14] Dan Hendrycks, Norman Mu, Ekin D Cubuk, Barret Zoph, Justin Gilmer, and Balaji Lakshminarayanan, "Augmix: A simple data processing method to improve robustness and uncertainty," *arXiv preprint arXiv:1912.02781*, 2019.

[15] Ekin D Cubuk, Barret Zoph, Jonathon Shlens, and Quoc V Le, "Randaugment: Practical automated data augmentation with a reduced search space," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, 2020, pp. 702–703.

[16] Piotr Dollár, Ron Appel, Serge Belongie, and Pietro Perona, "Fast feature pyramids for object detection," *IEEE transactions on pattern analysis and machine intelligence*, vol. 36, no. 8, pp. 1532–1545, 2014.

[17] Ahmet M Bagci, Rashid Ansari, and William D Reynolds, "Low-complexity implementation of non-subsampled directional filter banks using polyphase representations and generalized separable processing," in *2007 IEEE International Conference on Electro/Information Technology*. IEEE, 2007, pp. 422–427.

[18] Alican Bozkurt, Alexander Suhre, and A Enis Cetin, "Multiscale directional-filtering-based method for follicular lymphoma grading," *Signal, Image and Video Processing*, vol. 8, no. 1, pp. 63–70, 2014.

[19] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.

[20] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo, "Swin transformer: Hierarchical vision transformer using shifted windows," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 10012–10022.

[21] François Chollet, "Xception: Deep learning with depthwise separable convolutions," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1251–1258.

[22] Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam, "Mobilenets: Efficient convolutional neural networks for mobile vision applications," *arXiv preprint arXiv:1704.04861*, 2017.

[23] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen, "Mobilenetv2: Inverted residuals and linear bottlenecks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 4510–4520.