# Automatic diagnosis of knee osteoarthritis severity using Swin transformer

Aymen Sekhri
asekhri@inttic.dz
Laboratoire PRISME, université
d'Orléans
Orléans, France

Marouane Tliba
marouane.tliba@univ-orleans.fr
Laboratoire PRISME, université
d'Orléans
Orléans, France

Mohamed Amine Kerkouri
mohamed-amine.kerkouri@univ-orleans.fr
Laboratoire PRISME, université
d'Orléans
Orléans, France

Yassine Nasser
yassine.nasser@univ-orleans.fr
Laboratoire PRISME, université
d'Orléans
Orléans, France

Aladine Chetouani
aladine.chetouani@univ-orleans.fr
Laboratoire PRISME, université
d'Orléans
Orléans, France

Alessandro Bruno
alessandro.bruno@iulm.it
IULM AI Lab, IULM University
Milan, Italy

Rachid Jennane
rachid.jennane@univ-orleans.fr
IDP laboratory, université d'Orléans
Orleans, France

## ABSTRACT

Knee osteoarthritis (KOA) is a widespread condition that can cause chronic pain and stiffness in the knee joint. Early detection and diagnosis are crucial for successful clinical intervention and management to prevent severe complications, such as loss of mobility. In this paper, we propose an automated approach that employs the Swin Transformer to predict the severity of KOA. Our model uses publicly available radiographic datasets with Kellgren and Lawrence scores to enable early detection and severity assessment. To improve the accuracy of our model, we employ a multi-prediction head architecture that utilizes multi-layer perceptron classifiers. Additionally, we introduce a novel training approach that reduces the data drift between multiple datasets to ensure the generalization ability of the model. The results of our experiments demonstrate the effectiveness and feasibility of our approach in predicting KOA severity accurately.

## KEYWORDS

Medical imaging, Knee osteoarthritis, Vision transformers, Self-attention.
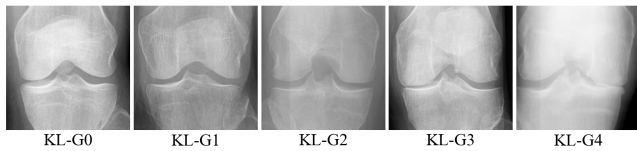
## 1 INTRODUCTION

Knee osteoarthritis (KOA) is a degenerative disease of the knee joint and the most common form of arthritis. It affects almost half of the population aged 65 years or older worldwide, causing pain, mobility limitation, and impaired quality of life. KOA is caused by a breakdown of knee articular cartilage and bone micro-architecture changes [7]. Joint space narrowing, osteophyte formation, and sclerosis are KOA's most visually relevant pathological features that can be visualized with radiographs. Although various imaging techniques such as magnetic resonance, computed tomography, and ultrasound have been introduced to diagnose osteoarthritis, radiography remains the most widely used method for initial diagnosis due to its accessibility, low cost, and widespread use.

Kellgren and Lawrence (KL) classified KOA severity into five stages based on the radiographic features, from KL-G0 for healthy cases to KL-G4 for severe cases [7] (See Fig 1). However, KOA changes gradually, so the evaluation into different stages is often subjective and depends on the operator. This causes subjectivity and makes the automatic KOA diagnosis a difficult task. In addition, the high similarity between the X-ray images increases the challenge of achieving an accurate diagnosis.

Several deep learning-based methods have been proposed for medical imaging applications [17], and many to diagnose KOA in recent years. In [1], Antony *et al.* employed Convolutional Neural Networks (CNNs) to quantify the severity of KOA from radiographic images. Their method is based on two main steps: first, automatically locate the knee joints using a Fully Convolutional Neural etwork (FCN), then, classify the knee joint images using a second CNN. In addition, to improve the quantification of KOA, they combined the classification loss with the regression loss to consider
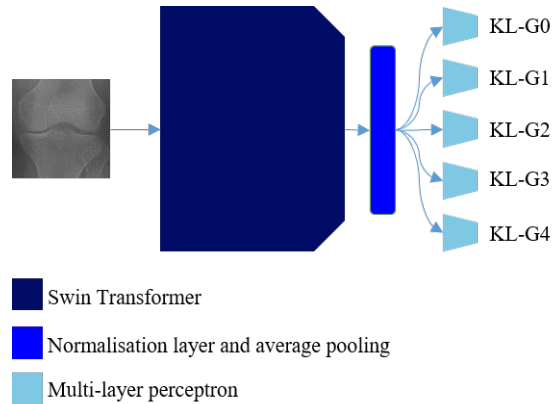
Figure 1: Samples of knee radiographs. KL-G0: healthy knee without osteoarthritis, KL-G1: doubtful osteoarthritis, KL-G2: minimal osteoarthritis, KL-G3: moderate osteoarthritis, and KL-G4: severe osteoarthritis.

the continuous aspect of the disease progression. Tuilpin *et al.* [16] presented a Siamese CNN network for KL grade prediction. They used three models with different random seeds and combined their outputs with a softmax layer to obtain the final KL grade. Chen *et al.* [5] proposed an ordinal loss for fine-tuning various CNN models to classify KOA severity. They leveraged the ordinal nature of the knee KL grading system and penalized incorrect classifications more by increasing the distance between the real and predicted KL grades. Nasser et al. [11] proposed a Discriminative Regularized Auto-Encoder (DRAE) for early KOA prediction using X-ray images. The proposed model uses a discriminative penalty term and the traditional AE reconstruction cost function to enhance the separability of the features learned from different classes. The aim was to boost the recognition system's performance by minimizing the inter-class variance and maximizing the intra-class distance. Recently, transformers have shown promising results in various medical imaging tasks [13]. Wang *et al.* [20] proposed a novel data augmentation method for early detection of KOA using a Vision Transformer model. The method involves shuffling the position embedding of non-ROI patches and exchanging the ROI patches with other images. The authors also used a hybrid loss function that combines label smoothing and cross-entropy to improve the model's generalization capability and avoid over-fitting. Several important studies [3],[6], [12], [14], [1], used two multi-center databases, the Osteoarthritis Initiative (OAI, https://nda.nih.gov/oai/) and the Multicenter Osteoarthritis Study (MOST, https://most.ucsf.edu/) by not accounting for the data drift problem. The latter occurs when a machine learning model trained on one dataset lowers its performance when tested on another set of data. Subsequently, data drift causes poor generalization and performance degradation.

In this work, we first investigate the use of the Swin transformer in predicting KOA severity from radiographic images. In particular, the Swin transformer is the core network that extracts high-level features and detects KOA-induced changes. Second, we introduce a multi-predictive classification header to address the high similarity problem between different KOA grades. In addition, to reduce the data drift problems between the data in the two databases, OAI and MOST, we tested several learning strategies to find the one providing the model with better generalization capabilities and balanced classification results.

The remainder of the paper is organized as follows: the proposed method is described in Section 2. Next, the obtained experimental results are presented in Section 3. Finally, the conclusions and outlooks are given in Section 4.



Figure 2: Swin transformer architecture with a multiple prediction head architecture

## 2 PROPOSED METHOD

The method proposed in this paper consists of two parts: 1) a Swin transformer as a features extractor and 2) a multi-prediction head network as a classifier. The schematic illustration of our proposed network is presented in Figure 2.

### 2.1 Swin Transformer

The Swin Transformer [9] is a state-of-the-art model that has been specifically designed to address the challenges of applying transformer models in the visual domain. While transformers have been widely successful in natural language processing, they have been less effective in computer vision due to the unique characteristics of visual data. The Swin Transformer proposes a novel architecture that leverages hierarchical feature maps and shift-based windows to improve the efficiency and performance of the model. With its innovative approach, the Swin Transformer has emerged as one of the most efficient and effective transformer models for visual applications. The model is divided into four stages, where the features are hierarchically extracted in each stage.

The input image with dimensions $H \times W \times 3$ is divided into $\frac{H}{4} \times \frac{W}{4}$ non-overlapping patches as tokens of size $4 \times 4 \times 3 = 48$. These tokens are then passed through the first stage, consisting of a linear embedding layer and two Swin Transformer blocks. The linear embedding layer projects the tokens into a higher-dimensional space denoted by $C$; after that, in the first Swin Transformer block, the multi-headed window self-attention mechanism (W-MSA) is employed. This mechanism computes self-attention only between patches within the same window, where each window contains $M \times M$ patches. The second Swin Transformer block utilizes shifted window multi-headed self-attention (SW-MSA), in which the partitioning windows are shifted by ($\lfloor \frac{M}{2} \rfloor$, $\lfloor \frac{M}{2} \rfloor$) patches with respect to the standard partitioning windows used in the previous block. This approach aims to create more relationships between neighboring patches previously located in different windows and reduce the computational complexity of the global MSA module used in vision transformer.

In the second stage, a patch merging layer is applied to group each $2 \times 2$ neighboring patches into a single patch of length $4C$,

thus reducing the number of patches to $\frac{H}{8} \times \frac{W}{8}$. These patches are then linearly projected to a dimension of size $2C$ and passed to two Swin Transformer blocks as in the first stage.

This process is repeated in the third stage, using 18 Swin Transformer blocks to produce $\frac{H}{16} \times \frac{W}{16}$ patches of length $4C$. Finally, in the fourth stage, two Swin Transformer blocks are used to produce $\frac{H}{32} \times \frac{W}{32}$ of length $8C$. These consecutive stages jointly produced a hierarchical representation like those of typical convolutional networks.

## 2.2 Multi-Prediction Head Network

The main task of our designed model is to be able to predict the KOA severity grade. This presents a case of a multi-class classification task. Traditionally this is solved by using a single MLP classification head with 5 outputs activated by a softmax function.

The complex nature of X-ray images imposes a high similarity between the images of adjacent KL Grades as shown in Figure 1. To address this issue, we decompose the task into multiple binary classification tasks. We use 5 MLP networks, each specializing in predicting one KL-Grade. This enhances the model's ability to extract and filter a rich representation for each class.

Let $f : X \rightarrow Z$ be our feature extractor, where $X$ and $Z$ are the input and latent spaces, respectively. $x$ represents the input image and $y$ their corresponding one hot encoding label. The predictive label $\hat{y}_i$ at the head classifier $MLP_i$ is defined as:

$$\hat{y}_i = MLP_i(f(x)) \tag{1}$$

The final predictive label $\hat{y}$ is computed then as follows:

$$\hat{y} = argmax(\bigcup_{i=0}^{4} \hat{y}_i) \tag{2}$$

where $i \in \{0 \dots 4\}$ represents the KL grades.

To sum up, our final model consists of a basic Swin-B encoder with $C = 128$ and $2, 2, 18, 2$ Swin Transformer blocks, followed by Normalisation and average pooling layers to produce a final representation vector of size 1024. This vector is then passed to 5 MLPs, one for each KL grade. Each MLP contains 3 linear layers of size 384, 48, 48, 1, respectively. The final layer of each MLP network has a single neuron to predict the occurrence probability of each grade.

## 2.3 Data Drift Correction

In this paper, we employ 2 of the most widely used datasets for KOA classification (i.e. MOST and OAI datasets). These datasets were collected over a substantial amount of time, from several medical centers, and were annotated by a multitude of medical practitioners. The inherent disparity of equipment, study subjects, radiography, and diagnostics methods between different medical centers caused a shift between the datasets as further discussed in Section 3.4.

We represent our model using the formula $h = g \circ f$, where $f : X \rightarrow Z$ and $g : Z \rightarrow Y$, represent the feature extractor and the multi-classification head, respectively. $X$ is the input image, $Z$ is the latent feature space, and $Y$ represents the label space.

To address the issue of data drift between the MOST and OAI datasets, we need to align the latent representational spaces between $Z_{MOST}$ and $Z_{OAI}$. This means that the feature extractor $f$ needs to be able to perceive the data distributions from $\mathcal{D}_{MOST}$ and $\mathcal{D}_{OAI}$ as belonging to the same distribution $\mathcal{D}$. It models relevant mutual features while discarding any dataset-specific information that could be considered noisy. This could be represented using the following equation:

$$\mathcal{D} = (\mathcal{D}_{MOST} \cup \mathcal{D}_{OAI}) \smallsetminus (\mathcal{N}_{MOST} \cup \mathcal{N}_{OAI}) \tag{3}$$

where $\mathcal{N}_{MOST}$ and $\mathcal{N}_{OAI}$ represent the noisy distribution of information specific to the MOST and OAI datasets, respectively.

To achieve this result, we train the model $h$ on the MOST dataset and then freeze the MLP layers $g$. We continue to train the feature extractor $f$ on the OAI dataset. This way, we force the feature extractor $f$ to align the representational space for both datasets. This proposed approach leverages the pre-trained source model effectively and adapts it to the target dataset by minimizing the shift between the data distributions in the latent representational space $Z$. The objective is to achieve this without compromising the prior knowledge of the pre-trained classifier.

## 2.4 Implementation

In order to train the model, we used the AdamW optimizer [10] with a learning rate of $3e - 5$, a weight decay of 0.05, an epsilon of $1e - 8$, and betas of $(0.9, 0.999)$ to adjust the weights. We trained the model with a batch size of 32 images for 300 epochs. We implemented the code in PyTorch and used an NVIDIA RTX A4000 GPU with 16 GB of VRAM to speed up the training process.

We also implemented various data augmentation techniques such as 15-degree rotation, translation, scaling, random horizontal flipping, and contrast adjustment with a factor of 0.3. These techniques have previously been used in similar studies to improve the performance of deep learning models on image classification tasks in order to address the problem of limited data and overfitting.

## 3 EXPERIMENTAL RESULTS

To evaluate the efficacy of the proposed approach, we conducted five experiments, described in this section.

### 3.1 Datasets

In this study, we employed two widely used and publicly available datasets:

**MOST dataset:** It contains 18,269 knee images that were segmented in the same manner as in [16]. We divided this dataset into three subsets, namely training, validation, and testing with a ratio of 6:1:3. Table 1 provides a summary of the dataset's partitioning. We use this dataset to train and evaluate our model's performance on knee image classification.

**OAI dataset:** It consists of 8260 already prepared knee images [5]. It is randomly divided into three subsets, namely training, validation, and testing with a ratio of 7:1:2. Table 2 summarizes the partitioning of the OAI dataset. We use this dataset to validate and test our model's performance.

| | KL-G0 | KL-G1 | KL-G2 | KL-G3 | KL-G4 | Total |
|---|---|---|---|---|---|---|
| Training | 4380 | 1759 | 1827 | 1986 | 1008 | **10960** |
| Validation | 730 | 294 | 304 | 331 | 168 | **1827** |
| Testing | 2190 | 880 | 914 | 994 | 504 | **5482** |

**Table 1: Label distribution of the MOST dataset**

| | KL-G0 | KL-G1 | KL-G2 | KL-G3 | KL-G4 | Total |
|---|---|---|---|---|---|---|
| Training | 2286 | 1046 | 1516 | 757 | 173 | **5778** |
| Validation | 328 | 153 | 212 | 106 | 27 | **826** |
| Testing | 639 | 296 | 447 | 223 | 51 | **1656** |

**Table 2: Label distribution of the OAI dataset**

## 3.2 Experimental Protocol

During the development of our model, we tested multiple configurations and compared them. In the first experiment, we use a single classifier to predict all grades simultaneously. In the second experiment, we use the same settings but employed the Multi-prediction head architecture, which involves breaking down the multi-classification problem into sub-binary classifications. For experiments three and four, we explored the data drift between two datasets by training only one dataset per experiment. Finally, in the fifth experiment, we tackled the issue of data drift by transferring the knowledge from the trained classifier on the source dataset (MOST) and solely training the feature extractor of our model on the target dataset (OAI).
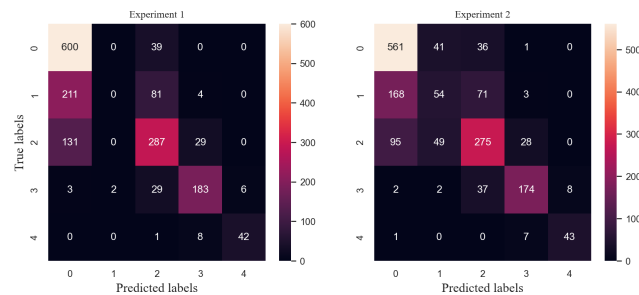
## 3.3 Quantitative Evaluation



**Figure 3: Confusion matrices on the OAI test set of Experiment 1 and 2.**

| Exp. | MOST | | OAI | |
|---|---|---|---|---|
| | Acc (%) ↑ | F1 ↑ | Acc (%) ↑ | F1 ↑ |
| 1 | 71.93 | 0.622 | 67.15 | 0.615 |
| 2 | 73.13 | 0.684 | 66.85 | 0.657 |
| 3 | 39.95 | 0.114 | 38.59 | 0.111 |
| 4 | **75.43** | **0.714** | 62.86 | 0.615 |
| 5 | 73.25 | 0.667 | **70.17** | **0.671** |

**Table 3: Comparison of the five experiments in terms of accuracy and F1 score on the OAI and MOST test sets.**

The performances obtained for each considered configurations are presented in Table 3. In the first two experiments, we observed an improvement in the F1 score for our model when using the Multi-prediction head architecture in the second experiment. Specifically, the model yielded a 0.062 and 0.042 F1 score increase compared to

the first experiment in the MOST and OAI test sets, respectively. We also notice an increase in accuracy on the MOST dataset.

Moreover, as seen by the confusion matrices in Figure 3, the architecture proposed in experiment 2 was able to avoid the catastrophic failure of detecting the KL-G1 observed in experiment 1. The grad KL-G1 is notoriously challenging to detect even for trained doctors due to the high similarity with the KL-G0 and KL-G2. In fact, the model correctly predicted 54 images in KL-G1 in experiment 2, while 0 images were classified in experiment 1. These results highlight the impact of dividing the multi-classification problem into sub-binary classification problems as described in sections 2.2. The substantial drop of performance in experiment 3 on both datasets is mainly attributed to the lack of a sufficient quantity of data. Transformer-based models are known to require a lot of data for training [4]. This has led to the underfitting of our model as it was not able to extract meaningful representations from this dataset. On the other hand, we notice that the performance of the model on the MOST dataset is quite similar, this is due to the richness of the representations in this dataset. In experiment 4, the MOST dataset contains more samples that cover a broader range of KOA severity levels than the OAI dataset as shown in Table 1. Consequently, MOST provides a more diverse and representative training set for our model, leading to better performance in the MOST test set. However, we still see a greater decrease in performance on the OAI dataset compared to experiment 2 in terms of accuracy and F1 score. Experiment 5 showed a considerable enhancement in performances on the OAI dataset compared to all other experiments, achieving a 70.17% accuracy and 0.671 F1-score, as shown in Table 3, while maintaining a high accuracy on the MOST dataset. This particularly highlights the significance and effectiveness of our method to reduce the data drift and align the latent representations of both datasets as described in section 3.4.

## 3.4 Latent Representation Ability

The reduction of the data drift is an important task for our model as shown in the previous quantitative results. Figure 4 depicts the distribution of latent features extracted for the samples of each dataset across the models produced through our previous experiments. We used the t-SNE algorithm [18] in order to reduce the dimensionality of the features. The data drift in the representation of the two datasets is clearly apparent for both experiments 1 and 2. Even though experiment 2 achieved better results, we still noticed the high disparity of performance between datasets. Due to the underfitting of the model in experiment 3, it was also unable to address the data drift. In experiment 4 the model was trained only on the MOST dataset. Because of the availability of data, we noticed a better general alignment for data distribution between datasets. But Figure 5 shows that the shift on the scale of individual classes is still noticeable. In experiment 5, we noticed a very strong alignment for both datasets on the general and class-specific levels in Figures 4 and 5, respectively. Our approach successfully aligned all the data points from both datasets, effectively mitigating the data drift problem. As a result, the learned representations were more relevant to the task, and the model's performance improved significantly.
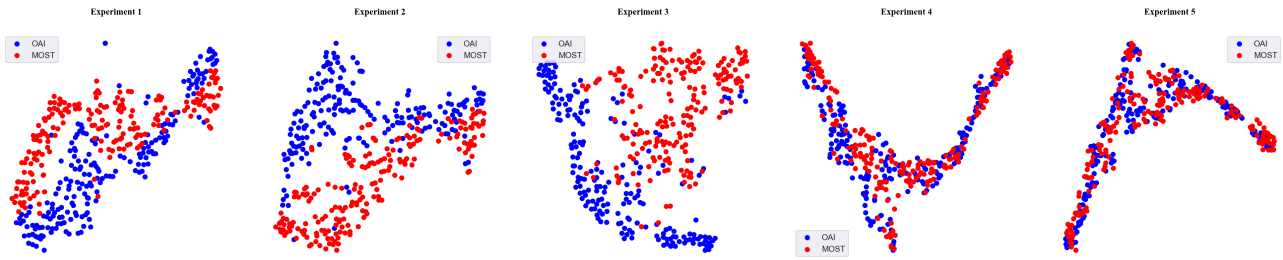
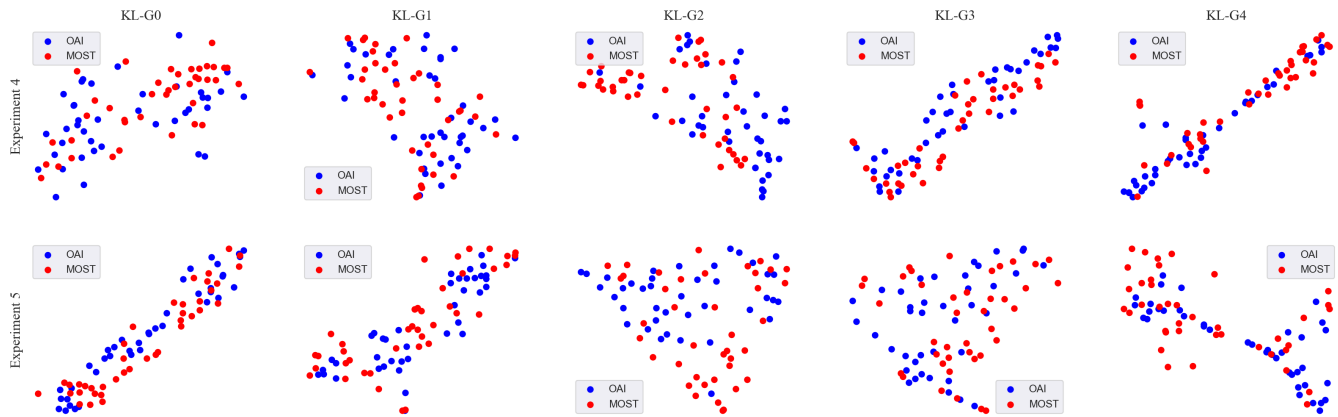**Figure 4: t-SNE visualizations of features learned by the model in each dataset.**



**Figure 5: t-SNE visualizations of features learned by the model in experiments 4 and 5 for each dataset.**
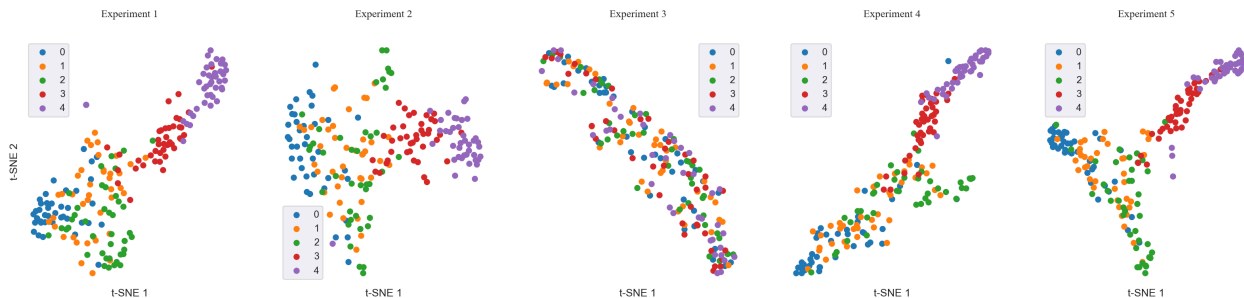


**Figure 6: t-SNE visualizations of grades separability in the OAI testset**

Figure 6 illustrates the distribution of latent representations of each class for each of our previous experiments on the OAI test-set. It highlights the ability of the model to discriminate and separate the different classes of KL-Grade. In experiment 3 where the underfitting occurred, we can observe the inability of the model to separate the distributions of the different classes. In experiments 1,2 and 4, the models were able to clearly separate the distributions of KL-G3 and KL-G4. Separating the KL-G0, KL-G1, and KL-G2 grades was more challenging in the first experiment due to the significant similarity between them and the use of a single MLP classifier. Along with the ability to align the distributions of both datasets, we noticed in Experiment 5 a better separability between

KL-G0, KL-G1, and KL-G2 which posed a challenge in other experiments. We observed a clear ability to discriminate between KL-G1 and KL-G2 especially, while KL-G0 and KL-G1 still pose some challenges because they represent the none existence and the very early stages of OA respectively.

Overall, these results demonstrate the effectiveness of our method in handling data drifts and enhancing the model's ability to differentiate between grades of KOA.

## 3.5 Qualitative Evaluation

We use GradCAM as a tool for interpretability purposes. By visualizing the last layer's activations of the feature extractor, we chose a sample from each grade, where the true labels of samples from (a)
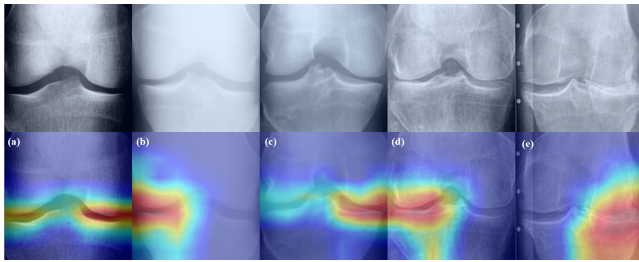
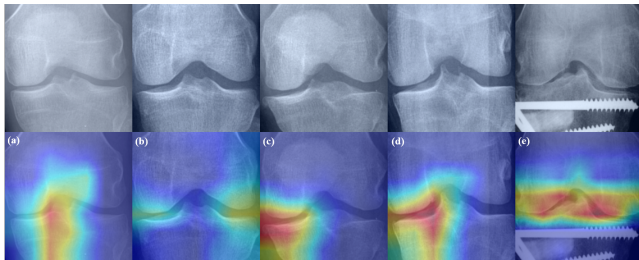**Figure 7: GradCAM of the corrected classified images.**



**Figure 8: GradCAM of the misclassifed images. (a) Predicted KL-G1 (b) Predicted KL-G0 (c) Predicted KL-G1 (d) Predicted KL-G2 (e) Predicted KL-G3.**

to (e) are from KL-G0 to KL-G4, respectively, as shown in Figures 7 and 8.

In Figure 7, we observed that the model effectively identified areas like osteophytes, joint space narrowing, and sclerosis, which are essential factors for assessing the severity of KOA [8]. This points out that our model bases its classifications on the right regions of interest commonly used in clinical diagnosis and not on non-relevant features.

Figure 8 represents misclassified samples. As can be observed, the model still focuses on the relevant regions around the knee joint. For instance, the model predicts sample (a) as KL-G1, even though the true KL grade was zero. It focused on the area where a medial joint space narrowing was present, which is a possible feature of KL-G1. Similar misclassifications occurred for samples (b), (c), and (d), where the model either overestimated or underestimated the KL grade, indicating the challenge of distinguishing between grades due to their high similarity and also the fact that the KL grade suffers from subjectivity/ambiguity among experts [15]. In sample (e), we encountered an image that contained an unusual object (i.e. A screw) in the tibia, which could potentially distract the model from the areas of the image that are crucial for grading KOA. However, our model demonstrated robustness by still being able to focus on the region of interest. Furthermore, our model classified the image as a KL-G3 instead of KL-G4, which are close compared to other KL-Grades. This result highlights the ability of our model to prioritize task-specific important features in the image and not be affected by irrelevant and noisy distractors.

## 3.6 State-of-the-art Comparison

Table 4 presents a comparison of the results obtained with state-of-the-art methods. We note that the methods used in these studies were trained differently. Specifically, some methods used the OAI training set exclusively, others used the MOST training set exclusively, and others used both bases. This diversity in learning can have an impact on the overall performance, and should therefore be carefully considered when interpreting the results.

Antony et al. [2] and [1] achieved accuracies of 53.40% and 63.60%, respectively, and F1-scores of 0.43 and 0.59, respectively. Chen et al. [5] used ordinal loss with different deep learning architectures and achieved accuracies of 69.60%, 66.20%, and 65.50% with Vgg19, ResNet50, and ResNet101, respectively, but they did not report F1-score. Tiulpin et al. [16] used a Siamese network and reported an accuracy of 66.71%. Wang et al. [19] achieved an accuracy of 69.18%.

Our proposed method, experiment 5, outperformed all other methods with an accuracy of 70.17% and an F1-score of 0.67. These results indicate the potential of our proposed method for improving the accuracy and reliability of knee osteoarthritis diagnosis, which could be valuable in clinical practice.

| Method | Acc (%) ↑ | F1 ↑ |
|---|---|---|
| Antony et al. 2016 [2] | 53.40 | 0.43 |
| Antony et al. 2017 [1] | 63.60 | 0.59 |
| Ordinal Loss (Vgg19) [5] | 69.60 | - |
| Ordinal Loss (ResNet50) [5] | 66.20 | - |
| Ordinal Loss (ResNet101) [5] | 65.50 | - |
| Siamese net [16] | 66.71 | - |
| Wang et al. [19] | 69.18 | - |
| **Ours experiment 5** | **70.17** | **0.67** |

**Table 4: Results for OAI dataset.**

## 4 CONCLUSION

In this paper, we proposed a new method to predict the severity of Knee OA from radiographic images using the Swin Transformer. Our results showed that this method achieved state-of-the-art performance on the OAI test set, significantly outperforming existing methods. We show that the Swin Transformer network is effective in extracting relevant knee OA information, which can be used to detect most of the symptoms of the disease. In addition, handling the data drift and using the multi-prediction head architecture significantly improves the accuracy of the model and helps reduce the similarity between features of nearby grades. Prospects for future work may involve other imaging modalities such as MRI, while exploring clinical and demographic data, to further improve the prediction of KOA severity.

## ACKNOWLEDGMENTS

## REFERENCES

[1] Joseph Antony, Kevin McGuinness, Kieran Moran, and Noel E O'Connor. 2017. Automatic detection of knee joints and quantification of knee osteoarthritis severity using convolutional neural networks. In *Machine Learning and Data Mining in Pattern Recognition: 13th International Conference, MLDM 2017, New York, NY, USA, July 15-20, 2017, Proceedings 13*. Springer, 376–390.

[2] Joseph Antony, Kevin McGuinness, Noel E O'Connor, and Kieran Moran. 2016. Quantifying radiographic knee osteoarthritis severity using deep convolutional neural networks. In *2016 23rd International Conference on Pattern Recognition (ICPR)*. IEEE, 1195–1200.

[3] Abdelbasset Brahim, Rachid Jennane, Rabia Riad, Thomas Janvier, Laila Khedher, Hechmi Toumi, and Eric Lespessailles. 2019. A decision support tool for early detection of knee OsteoArthritis using X-ray imaging and machine learning: Data from the OsteoArthritis Initiative. *Computerized Medical Imaging and Graphics* 73 (2019), 11–18.

[4] Yun-Hao Cao, Hao Yu, and Jianxin Wu. 2022. Training vision transformers with only 2040 images. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXV*. Springer, 220–237.

[5] Pingjun Chen, Linlin Gao, Xiaoshuang Shi, Kyle Allen, and Lin Yang. 2019. Fully automatic knee osteoarthritis severity grading using deep neural networks with a novel ordinal loss. *Computerized Medical Imaging and Graphics* 75 (2019), 84–92.

[6] AA Gatti. 2018. NEURALSEG: state-of-the-art cartilage segmentation using deep learning–analyses of data from the osteoarthritis initiative. *Osteoarthritis and Cartilage* 26 (2018), S47–S48.

[7] Mark D Kohn, Adam A Sassoon, and Navin D Fernando. 2016. Classifications in brief: Kellgren-Lawrence classification of osteoarthritis. *Clinical Orthopaedics and Related Research*® 474 (2016), 1886–1893.

[8] Michelle J Lespasio, Nicolas S Piuzzi, M Elaine Husni, George F Muschler, AJ Guarino, and Michael A Mont. 2017. Knee osteoarthritis: a primer. *The Permanente Journal* 21 (2017).

[9] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. 2021. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF international conference on computer vision*. 10012–10022.

[10] Ilya Loshchilov and Frank Hutter. 2017. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101* (2017).

[11] Yassine Nasser, Rachid Jennane, Aladine Chetouani, Eric Lespessailles, and Mohammed El Hassouni. 2020. Discriminative Regularized Auto-Encoder for early

[12] Rabia Riad, Rachid Jennane, Abdelbasset Brahim, Thomas Janvier, Hechmi Toumi, and Eric Lespessailles. 2018. Texture analysis using complex wavelet decomposition for knee osteoarthritis detection: Data from the osteoarthritis initiative. *Computers & Electrical Engineering* 68 (2018), 181–191.

[13] Fahad Shamshad, Salman Khan, Syed Waqas Zamir, Muhammad Haris Khan, Munawar Hayat, Fahad Shahbaz Khan, and Huazhu Fu. 2023. Transformers in medical imaging: A survey. *Medical Image Analysis* (2023), 102802.

[14] Albert Swiecicki, Nianyi Li, Jonathan O'Donnell, Nicholas Said, Jichen Yang, Richard C Mather, William A Jiranek, and Maciej A Mazurowski. 2021. Deep learning-based algorithm for assessment of knee osteoarthritis severity in radiographs matches performance of radiologists. *Computers in biology and medicine* 133 (2021), 104334.

[15] Aleksei Tiulpin and Simo Saarakkala. 2020. Automatic grading of individual knee osteoarthritis features in plain radiographs using deep convolutional neural networks. *Diagnostics* 10, 11 (2020), 932.

[16] Aleksei Tiulpin, Jérôme Thevenot, Esa Rahtu, Petri Lehenkari, and Simo Saarakkala. 2018. Automatic knee osteoarthritis diagnosis from plain radiographs: a deep learning-based approach. *Scientific reports* 8, 1 (2018), 1–10.

[17] Marouane Tliba, Aymen Sekhri, Mohamed Amine Kerkouri, and Aladine Chetouani. 2022. Deep-Based Quality Assessment of Medical Images Through Domain Adaptation. In *2022 IEEE International Conference on Image Processing (ICIP)*. 3692–3696. https://doi.org/10.1109/ICIP46576.2022.9897600

[18] Laurens Van der Maaten and Geoffrey Hinton. 2008. Visualizing data using t-SNE. *Journal of machine learning research* 9, 11 (2008).

[19] Yifan Wang, Xianan Wang, Tianning Gao, Le Du, and Wei Liu. 2021. An automatic knee osteoarthritis diagnosis method based on deep learning: data from the osteoarthritis initiative. *Journal of Healthcare Engineering* 2021 (2021), 1–10.

[20] Zhe Wang, Aladine Chetouani, and Rachid Jennane. 2023. Transformer with Selective Shuffled Position Embedding using ROI-Exchange Strategy for Early Detection of Knee Osteoarthritis. *arXiv preprint arXiv:2304.08364* (2023).