

1 **ARTIFICIAL INTELLIGENCE ACCURATELY DETECTS TRAUMATIC**  
2 **THORACOLUMBAR FRACTURES ON SAGITTAL RADIOGRAPHS**

3 **AUTHORS**

4 Guillermo Sanchez Rosenberg<sup>1,5</sup>, MD ORCID iD: 0000-0003-2304-7568

5 Andrea Cina<sup>2</sup> ORCID iD: 0000-0002-6016-6200

6 Giuseppe Rosario Schirò<sup>3</sup>, MD ORCID iD: 0000-0001-5190-5131

7 Pietro Domenico Giorgi<sup>3</sup>, MD ORCID iD: 0000-0002-3682-112

8 Prof. Boyko Gueorguiev<sup>1</sup> ORCID iD: 0000-0001-9795-115X

9 Prof. Mauro Alini<sup>1</sup> ORCID iD: 0000-0002-0262-1412

10 Peter Varga<sup>1</sup>, PhD ORCID iD 0000-0003-2738-6436

11 Fabio Galbusera<sup>2</sup>, PhD ORCID iD: 0000-0003-1826-9190

12 Enrico Gallazzi<sup>4</sup>, MD ORCID iD: 0000-0001-9287-9937

13 **AFFILIATIONS**

14 1. AO Research Institute, Davos, Switzerland

15 2. IRCCS Istituto Ortopedico Galeazzi, Milano, Italy

16 3. ASST GOM Niguarda, Milano, Italy

17 4. Centro Specialistico Ortopedico Traumatologico G. Pini – CTO, Milano, Italy

18 5. Spital Lachen, Lachen, Switzerland

19 **CORRESPONDING AUTHOR**

20 Guillermo Sánchez Rosenberg, M.D.

21 [gsanchez@ufm.edu](mailto:gsanchez@ufm.edu)

22 Research Fellow

23 AO Research Institute Davos, Switzerland

24 Present Address;

25 Chirurgie, Spital Lachen

26 Oberdorfstrasse 41

27 8853, Lachen

28 Schwyz, Switzerland

29 Financial Disclosure: This study was performed with the assistance of the AO Foundation.

30

1 **Abstract**

2

3 **Background context.** Traumatic thoracolumbar (TL) fractures are frequently encountered in  
4 emergency rooms. Sagittal and anteroposterior radiographs are the first step in the trauma routine  
5 imaging. Up to 30% of TL fractures are missed in this imaging modality, thus requiring a CT and/or  
6 MRI to confirm the diagnosis. A delay in treatment leads to increased morbidity, mortality,  
7 exposure to ionizing radiation and financial burden. Fracture detection with Machine Learning  
8 models has achieved expert level performance in previous studies. Reliably detecting vertebral  
9 fractures in simple radiographic projections would have a significant clinical and financial impact.

10 **Purpose.** To develop a deep learning model that detects traumatic fractures on sagittal  
11 radiographs of the TL spine.

12 **Study design/setting.** Retrospective Cohort study.

13 **Methods.** We collected sagittal radiographs, CT and MRI scans of the TL spine of 362 patients  
14 exhibiting traumatic vertebral fractures. Cases were excluded when CT and/or MRI were not  
15 available. The reference standard was set by an expert group of three spine surgeons who  
16 conjointly annotated the sagittal radiographs of 171 cases. CT and/or MRI were reviewed to  
17 confirm the presence and type of the fracture in all cases. 302 cropped vertebral images were  
18 labelled 'fracture' and 328 'no fracture'. After augmentation, this dataset was then used to train,  
19 validate, and test deep learning classifiers based on ResNet18 and VGG16 architectures. To  
20 ensure that the model's prediction was based on the correct identification of the fracture zone, an  
21 Activation Map analysis was conducted.

22 **Results:** Vertebrae T12 to L2 were the most frequently involved, accounting for 48% of the  
23 fractures. A4, A3 and A1 were the most frequent AO Spine fracture types. Accuracies of 88% and  
24 84% were obtained with ResNet18 and VGG16 respectively. The sensitivity was 89% with both  
25 architectures but ResNet18 showed a higher specificity (88%) compared to VGG16 (79%). The  
26 fracture zone was precisely identified in 81% of the heatmaps.

1 **Conclusions.** Our AI model can accurately identify anomalies suggestive of vertebral fractures in  
2 sagittal radiographs by precisely identifying the fracture zone within the vertebral body.

3 **Clinical significance.** Clinical implementation of a diagnosis aid tool specifically trained for TL  
4 fracture identification is anticipated to reduce the rate of missed vertebral fractures in emergency  
5 rooms.

6

## 7 **Keywords**

8 Thoracolumbar fractures; spinal fracture; vertebral fracture; bone fracture detection; heatmap;  
9 machine learning; deep learning; artificial intelligence

10

## 11 **Introduction**

12 The thoracolumbar (TL) spine is one of the most common site of traumatic fracture occurrence,  
13 with an incidence that ranges from 32 to 64/100.000 per year; furthermore, traumatic TL fractures  
14 have a rate of associated neurological injuries from 22% to 51% depending on the fracture type,  
15 and a require surgical treatment in 38% of the cases. (1–3) Traumatic TL fractures are associated  
16 with decreased physical function, severe reduction of quality of life and the lowest rate of return to  
17 work among all major organ injuries. (4) Additionally, the overall mortality associated to spinal  
18 injuries is 17%. (5) Besides this elevated disease burden and prevalence, the treatment of  
19 vertebral fractures is costly. The annual estimated economic cost, in the United States alone,  
20 surpassed the billion dollar figure already in 2011. (6)

21 The severity of traumatic TL fractures can range from a simple apophyseal fracture without  
22 structural impairment to a complete dislocation of the spine. Sagittal and anteroposterior  
23 radiographs are the first step in the trauma routine imaging. However, they are not a very reliable  
24 diagnosis aid when suspecting TL fractures: the worldwide reported false-negative rate is as high  
25 as 30%. (7) Moreover, the current classification systems such as the AOSpine Classification and  
26 the Thoracolumbar Injury Classification and Severity score (TLICS) stratify fracture severity on

1 parameters such as anterior failure of the vertebral disc under compression and posterior integrity  
2 of ligamentous structures (8), which are not readily identifiable in radiographs. The severity and  
3 instability of the fracture will determine the need for conservative or surgical treatment. Thus,  
4 surgeons must resort to second level imaging such as Computer Tomography (CT) or Magnetic  
5 Resonance Imaging (MRI) to determine the treatment strategy. (9) The need for second level  
6 imaging inevitably leads to a delay in diagnosis, ranging from hours to even months, frequently  
7 resulting in poorer clinical outcomes. (10,11)

8 The recent explosion of labeled data, namely 'big data', has brought upon the era of artificial  
9 intelligence (AI). Within the healthcare sector, AI is being now used for several applications  
10 including drug discovery, remote patient monitoring, risk management, wearables, virtual  
11 assistants, and hospital management. Regarding medical diagnostics and imaging, the field of  
12 radiology has been particularly benefited. (12) The management of patients with musculoskeletal  
13 diseases could be improved by these innovations, provided that optimal accuracies are preserved.  
14 (13–15) By supporting the treating physician in identifying anomalies on patients imaging studies,  
15 AI is posed to considerably reduce diagnostic errors (16), therefore improving clinical outcomes in  
16 the treatment of vertebral fractures.

17 Deep learning (DL) is a supervised machine learning method that uses an algorithmic structure  
18 based on neural networks, such as Convolutional Neural Networks (CNN) (17). This method has  
19 been reported to perform as good or even better than humans in image classification (18). The  
20 power of this technique lies in the ability to identify and extract relevant features from labeled data  
21 at a grand scale. (19)

22 Recently, several proof of concept papers were published showing the application of AI and DL in  
23 spine imaging, showing promising results in the evaluation of degenerative disorders (20), adult  
24 deformities (21), adolescent idiopathic scoliosis (22), detection of primary and secondary bone  
25 tumors (23,24), and vertebral fractures (25).

26 Implementing a diagnosis aid tool specially trained for fracture identification in the clinical practice  
27 is anticipated to reduce the rate of missed vertebral fractures in emergency rooms. Murata et al

1 have recently published a model capable of detecting vertebral fractures in plain radiographs. (26)  
2 However, Murata's study is limited to radiographs displaying only one vertebral fracture, therefore  
3 invalidating its application in a scenario where multiple are present. Of note, the annotation of each  
4 image subgroup was done by a single spine surgeon, which could result in the introduction of  
5 confirmation bias in the standard of reference.

6 The main aim of this study was to develop an AI-based algorithm capable of accurately detecting  
7 vertebral fractures in sagittal radiographs of the TL spine. The secondary aim was to gain a deeper  
8 understanding of the model's interpretation of the 'fracture zone' through a heatmap  
9 representation.

10

## 11 **Methods**

12 This work utilized a CNN-based supervised learning approach.(17) The model was trained on  
13 cropped single vertebrae. The standard of reference was set by three expert spine surgeons, who  
14 evaluated the plain sagittal radiograph together with second level imaging, namely CT and/or MRI  
15 and annotated the single vertebra image as 'fracture' or 'non-fracture'.

### 16 *Image acquisition and standard of reference*

17 Sagittal radiographs, CT and MRI scans from the TL spine of 362 patients of more than 12 years of  
18 age and exhibiting traumatic vertebral fractures were retrospectively collected from a Spine  
19 Surgery reference Center (ASST Grande Ospedale Metropolitano Niguarda, Milano, Italy). The  
20 Ethical Committee approval was granted via the *Comitato Etico Milano Area 3* under the Ref. No.  
21 359-24062020. To ensure accuracy of the diagnoses, only cases with an initial radiograph and a  
22 subsequent CT or MRI were included. Fractures resulting from mechanisms other than trauma,  
23 such as osteoporosis or pathologic fractures were excluded. Sagittal radiographs and second level  
24 imaging of 151 patients were selected for the final analysis. An expert group of three spine  
25 surgeons with more than 30 years of accumulated experience reviewed each case individually to  
26 identify the presence of a fracture, and then classified them according to the AO Spine

1 Classification. Finally, in the cases where disagreement existed, meetings were held to reach  
2 unanimous consensus for each case.

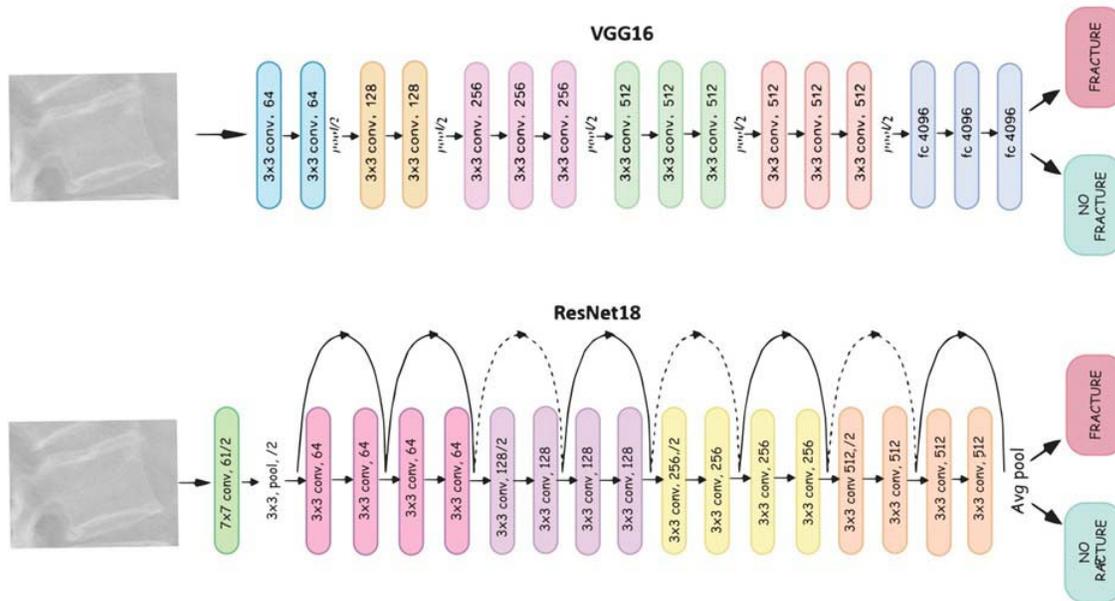
3 A total of 630 single vertebra images, obtained from in 222 sagittal radiographs of the TL spine  
4 were annotated. Of these, 302 were annotated as 'fracture' and 328 as 'non-fracture'. The  
5 annotation process was performed using a C++ software specifically developed for this project.  
6 The annotator indicated the Region of Interest (ROI) with a bounding box around a vertebra, and  
7 assigned the class, 'fracture' or 'non-fracture', the level and the corresponding AO Spine  
8 Classification.

#### 9 *Training and Test Sets:*

10 The image dataset was split into a training set (N = 578) and a test set (N = 52). To increase the  
11 generalization capability of the model, we used augmentation techniques such as random rotation,  
12 flipping, and shifting. In this way the model was trained on different versions of the same image  
13 during the training epochs.

#### 14 *Classification with Deep Learning*

15 We compared the performances of VGG16 (27) and ResNet18 (28) DL architectures, which were  
16 in the top three of the ImageNet challenge (29) in 2014 and 2015 respectively. These CNNs  
17 achieved state-of-the-art performances on computer vision tasks such as image classification,  
18 object detection, and landmark localization. (25) The difference between these CNNs is that  
19 VGG16 is a plain neural network with a deep sequence of convolutional layers followed by max-  
20 pooling, while ResNet functions with the so-called residual blocks, where the input of each block is  
21 summed to the output of the same block creating a skip connection (**Fig. 1**). The motivation behind  
22 the skipping connection is the vanishing gradient problem that arises during the training of very  
23 deep CNNs. (27) The images were resized to 512x512 pixels and normalized to have zero mean  
24 and unit variance, according to the image guidelines used in the ImageNet challenge.



1

2 **Figure 1:** DL network architectures used in this study. The VGG16 is a sequence of convolution  
3 and max-pooling operations. The number of parameters to learn is remarkably high (about 138  
4 million). ResNet18 presents the so-called Residual Blocks that represent the blocks of convolution  
5 operation between the skipping connections represented by the arrows, reducing the number of  
6 parameters to approximately 11 million.

7

8 Since the number of available images was not high, we used a technique called Transfer Learning  
9 (30,31) that consists of exploiting parameters of the 2 models that obtained state-of-the-art  
10 performances on the ImageNet challenge and retrained only the last few layers, the 2 last residual  
11 blocks in the present study, on the new task of vertebral fracture classification. We replaced the  
12 last fully connected layers of both networks (with 1000 neurons each) with a fully connected layer  
13 with 2 neurons representing the 2 classes. For the fully connected layer we used the softmax  
14 activation function that returns a probability distribution that assigns to each sample a probability of  
15 belonging to a class. The objective function is the negative log-likelihood loss used together with  
16 the softmax.

1 To find the best hyperparameters to train the network, in particular the learning rate and the batch  
2 size, we performed cross-validation with 10 folds where the training set is split into 10 parts, 9 of  
3 which are used to train the model and 1 for validation. This process was repeated 10 times  
4 iteratively until all the combinations of the folds have been used for the training/validation process.  
5 Finally, the entire training set was used for training with the best hyperparameters found in the  
6 cross-validation and we evaluated the performances of the model on the test set.

7 The model was implemented in Python language using PyTorch (34), a deep learning framework  
8 developed by Facebook. For the training and the evaluation, we used a Linux workstation with a  
9 NVIDIA QUADRO RTX 5000. The models ran for 200 epochs using a batch size of 32 and a  
10 learning rate of 0.00016, which resulted as the best hyperparameters in the cross-validation step.  
11 We used the Adam optimizer for model optimization and a method that reduced the learning rate  
12 by a factor of 0.1 if the accuracy did not improve for 10 epochs in a row (ReduceLROnPlateau in  
13 PyTorch).

#### 14 *Evaluation*

##### 15 Model's Performance Parameters

16 The models' performance was assessed quantitatively by calculating the accuracy, sensitivity, and  
17 specificity in fracture identification. Accuracy represents the ability of the model to assign the  
18 images to the correct class, described by the AUC (**Fig. 5**); the AUC outputs a value that displays  
19 the probability of a random sample being correctly classified by the algorithm, thus indicating the  
20 capacity of a classifier to distinguish between two classes. (32) Sensitivity describes the ability to  
21 detect the fractures, and specificity is the ability to detect the lack of a fracture.

##### 22 Understanding the Model's Prediction

23 To ensure that the model's prediction was based on the correct identification of the fracture zone,  
24 we implemented Activation Maps on the single vertebral images to highlight the image regions that  
25 lead the model to classify an image into a specific class. The activation maps were obtained by  
26 multiplying the second last layer of the neural network (the last feature maps) by the weights that  
27 point to the neuron of the class predicted by the model. This way, all pixels of the feature maps

1 were weighted according to the model prediction highlighting the most important parts of the  
2 images that determined the prediction. The heatmaps were evaluated by the same surgeons that  
3 set the standard of reference, judging whether the “warm zones” seen in the Activation Maps  
4 correlated to the fracture zones seen the CT or MRI images.

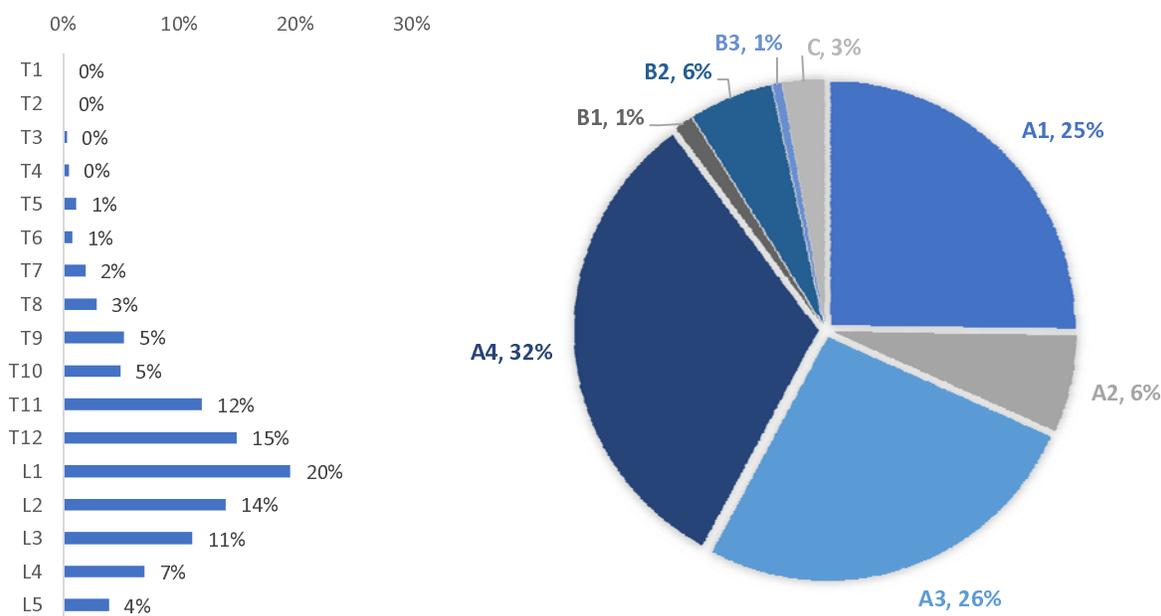
5

## 6 Results

### 7 Clinical Dataset

8 For the final analysis, a total of 151 cases of patients with TL fracture with availability of the initial  
9 radiograph and subsequent CT or MRI were selected. A total of 222 TL sagittal radiographs were  
10 analyzed and classified by the expert group of spine surgeons. Vertebrae T12 to L2 were the most  
11 frequently involved, accounting for 48% of the fractures (**Fig. 2**). Axial compression fractures,  
12 namely the AO Spine types A4, A3 and A1, were the most frequent injury mechanisms (**Fig. 3**).

13



14

15 **Figure 2:** Epidemiological distribution of TL  
16 fractures.

17 **Figure 3:** TL fracture type per AO  
18 Classification.

1

## 2 *Deep Learning Model*

3 Both DL architectures achieved high accuracy, sensitivity, and specificity after hyperparameter  
4 optimization, but ResNet18 performed better in all these aspects compared to VGG16. Both  
5 models predicted three false negatives (5.8%) by misclassifying three 'fracture' images as 'no  
6 fracture'. ResNet18 showed increased specificity, predicting three false positives in comparison to  
7 the 5 of the VGG16, namely classifying 'no fracture' images as 'fracture'. (**Fig. 4**).

8

		ResNet18		VGG16	
Ground Truth	no fracture	TN 25 48.1%	FP 3 5.8%	TN 25 48.1%	FP 5 9.6%
	fracture	FN 3 5.8%	TP 21 40.4%	FN 3 5.8%	TP 19 36.5%
		no fracture	fracture	no fracture	fracture
		Model		Model	

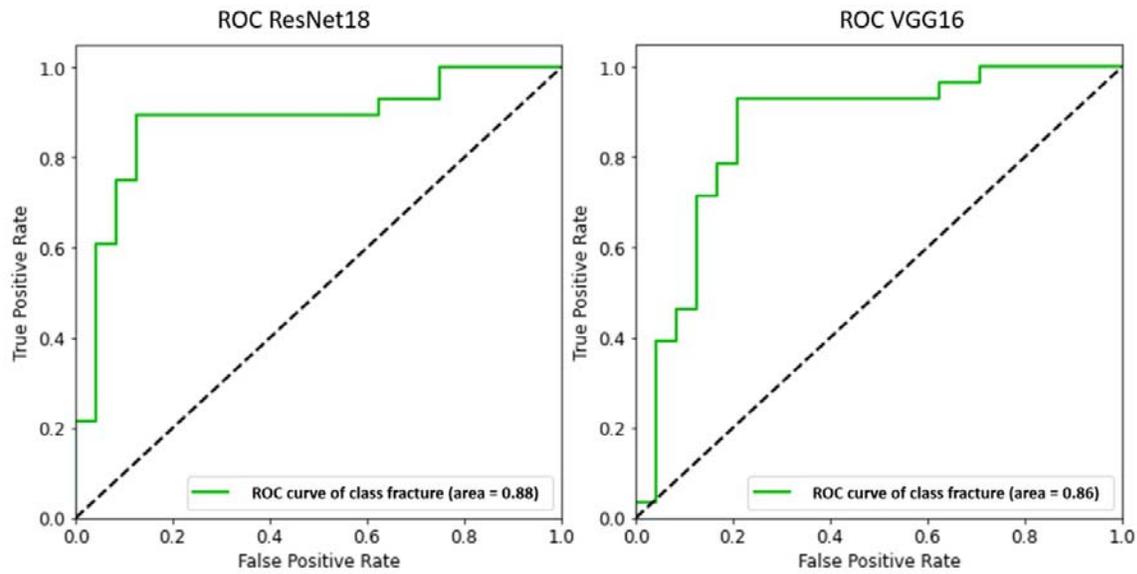
9

10 **Figure 4:** Confusion matrices obtained with the two DL architectures. The ResNet18 model made  
11 6 misclassifications, whereas VGG16 made 8. TN: True negative; FN: false negative; TP: true  
12 positive; FP; false positive

13

14 In terms of area under the ROC (**Fig. 5**) the ResNet18 performed better than the VGG16, with 0.88  
15 and 0.86 respectively for both classes.

16

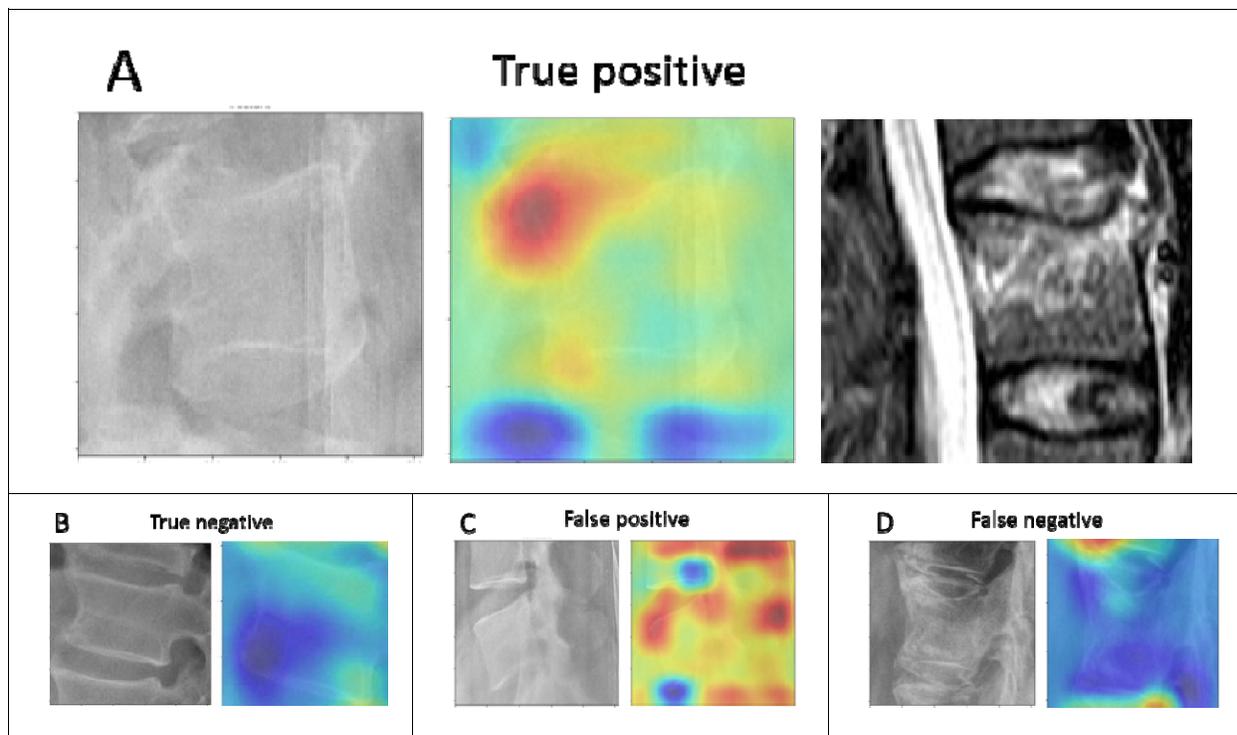


2 **Figure 5:** Comparison of the ROC curve obtained with ResNet18 and VGG16.

3

4 To ensure that the model's prediction was based on the correct identification of the fracture zone,  
5 we conducted an Activation Map analysis. The resulting heatmaps depict which areas of the  
6 images led the model to classify the vertebrae as 'fracture' or 'no fracture' (**Fig. 6**). In 81% of the  
7 single vertebrae, the "warm zone" correlated to the fracture zone observed in the corresponding  
8 CT or MRI. Interestingly, in two occasions, the model's prediction made the surgeons question the  
9 correctness of the ground truth. After verification via MRI and CT, the two images had to be  
10 reassigned to the opposite class. The model's prediction effectively amended human errors made  
11 during the annotation process. Accounting for this reassignment, the number of false negatives  
12 would be reduced from 3 to 1, thus increasing the sensitivity from 89% to 96%.

13



1

2 **Figure 6:** Heatmap analysis of the fracture zone. In Panel a, the heatmap correlates to the fracture  
3 zone identified in MRI. Panel b shows a true negative where the heatmap did not highlight a critical  
4 zone on the vertebra. Panel c incorrectly indicates presence of a fracture. Panel d was originally  
5 classified as 'fracture' and thus accounted for as a false negative, but it was then reclassified as  
6 'no-fracture' after the evaluation of the heatmap.

7

## 8 **Discussion**

9 This study demonstrated that AI-based techniques can detect vertebral fractures on radiographs  
10 with very high accuracy. Both models achieved similar sensitivity and specificity to that achieved by  
11 expert surgeons and radiologists (26,33–37) and the average sensitivity reported in a recent  
12 review.(38) ResNet18 showed better performance in identifying vertebral fractures compared to  
13 VGG16. Additionally, ResNet18 was less resource intensive in terms of memory used by the  
14 network parameters (43 MB vs 524 MB) and faster in the inference. To our knowledge, this is the  
15 first time that ResNet18 has been adapted for fracture identification purposes. Furthermore,

1 ResNet18 predictions were based in most cases, on the regions of vertebrae corresponding to the  
2 fracture zone observed in the CT and MRI. To ensure the quality of the diagnostic trial, the  
3 reference standard was established only after confirmation on CT and/or MRI by three different  
4 experts, as suggested by a recent meta-analysis of the diagnostic accuracy of deep learning in  
5 orthopedic fractures (38) and the expert panel recommendations from the Radiological Society of  
6 North America. (39)

### 7 *Comparison of Model's Performance*

8 Only recently have AI-based models been used to attempt fracture detection on radiographs.  
9 Presumably the first proof-of-concept paper using CNNs for fracture identification was published by  
10 Olczak in 2017. They compared the fracture identification capability of 5 existing CNNs. Fracture  
11 presence was deduced by extracting a combination of expressions and keywords from the  
12 radiologist's report, namely metadata. Contrasting to our results, VGG16 showed the best  
13 performance in their study, achieving 83% fracture identification accuracy. (33)

14 Other authors have also adapted CNNs to the problem of fracture detection, focusing exclusively  
15 on image interpretation, namely the information within the image file. Kim and MacKinnon used an  
16 adapted version of the Inception V3 model to identify distal radius fractures on sagittal radiographs,  
17 achieving an AUC of 0.954. (36) Although a significant limitation of this previous study was the  
18 exclusion of radiographs if the single lateral projection was inconclusive for the presence of  
19 fracture, their model analyzed the complete radiograph image instead of a cropped region of  
20 interest, as we and most other researchers have done. Chung et al used an adapted ResNet-152  
21 on cropped anteroposterior radiographs of the shoulder to distinguish fractured from normal  
22 humeri, achieving an accuracy of 95%, AUC of 0.996, sensitivity of 99% and specificity of 97% in  
23 the optimal cutoff point. (40) Adams et al. used cropped radiographs of surgically confirmed  
24 femoral neck fractures to compare the performances between AlexNet and GoogLeNet.  
25 GoogLeNet outperformed AlexNet, achieving an overall accuracy of 90.6%. Given that the  
26 reference standard was established surgically, the probability of bias introduction into the model  
27 was cleverly minimized. (34) Similarly, we minimized annotation bias by training the model

1 exclusively with radiographs where the presence of the fracture was confirmed via CT or MRI.  
2 Urakawa et al also evaluated cropped radiographs of femoral neck fractures and achieved an  
3 accuracy of 95.5%, AUC 0.984, sensitivity of 93.9% and specificity of 97.4% using an adapted  
4 version of VGG16. (37)

5 Recently, a model based on Visual Recognition V3 (IBM, Armonk, NY, USA) was used to identify  
6 vertebral fractures by Murata et al, achieving an accuracy, sensitivity, and specificity of 86.0%,  
7 84.7%, and 87.3% respectively. (26) While their results are similar to ours, there are important  
8 methodological differences to consider. To avoid introduction of systematic errors while training the  
9 model, all the fractures included in our study were evaluated individually by expert spine surgeons  
10 before annotation, and then discussed in consensus meetings where discrepancy occurred. In  
11 contrast, each classifying surgeon in the study of Murata et al. seemingly evaluated a single  
12 subgroup of images. While our model was trained to identify anomalies in single vertebrae to  
13 eliminate confounding factors and ensure a future clinical applicability, as shown in the heatmap  
14 analyses (**Fig. 5**), Murata's group analyzed the entire radiograph. The exclusion of cases with  
15 multiple traumatic fractures impairs the application of their model in the clinical practice. However,  
16 the inclusion of anteroposterior radiographs approaches a regular clinical scenario where both  
17 projections would be evaluated. In addition to the use of a different model, these factors might  
18 have contributed to the marginally better performance achieved in our study.

### 19 *Heatmap Analysis*

20 In 81% of the cases, our model's prediction of 'fracture' or 'no fracture' was based on a precise  
21 identification of the anomalies in single vertebrae, confirmed by correlating the "warmer zones" with  
22 the findings in CT and/or MRI (**Fig. 6**). Interestingly, two images which the model predicted as 'no  
23 fracture' by the model were originally classified as 'fracture' by the surgeons. A reassessment of  
24 the images supported the model's prediction and increased the model's sensitivity to 96%.  
25 Although this finding should be cautiously considered due to its exemplary nature, it illustrates the  
26 potential of AI to contribute to physicians' decisions in the clinical workflow.

## 1 *Limitations*

2 The present study has limitations. First, the dataset had a relatively small size. Traumatic vertebral  
3 fractures are commonly diagnosed based on CT or MRI only, obviating the need for radiographs in  
4 most cases. Although a larger database would arguably have enhanced the performance of the  
5 CNNs, we mitigated the impact of this limitation by performing aggressive image augmentation and  
6 taking advantage of models pre-trained on the ImageNet dataset. Since in the clinical workflow, a  
7 surgeon or physician would mostly rely on CT and MRI to confirm the presence of a fracture, a  
8 comparison with a model trained only with sagittal radiographs seemed unbalanced for this study's  
9 purpose. A comparative evaluation of the performance of the classifier and that of surgeons who  
10 are naïve to the clinical images will be reported in a future study. Regarding the heatmaps, it  
11 should be noted that the activation maps do not necessarily show the fracture zone but rather the  
12 zones that are more important in determining the output of the classifier, which may not correspond  
13 to the fracture itself.

## 14 *Clinical Relevance of AI for automated traumatic lesion detection*

15 Introducing systems of radiograph interpretation can reduce the frequency of misdiagnosis to  
16 below 0.3%.<sup>(41)</sup> Failures in fracture identification can be considerably reduced by implementing a  
17 second-stage verification algorithm at the end of the normal workflow to complement the  
18 interpretation of the physician. This way, introduction of bias or distractions would be avoided.

19 Contrary to common belief, computer-aided diagnostic tools are not necessarily aimed to replace  
20 expert human interpretation of medical imaging. In the authors view, the goal is minimization of  
21 Diagnostic Errors. Currently up to 30% vertebral fractures in radiographs are missed (1–3),  
22 resulting in either delayed or missed diagnosis. Both outcomes are qualified as Diagnostic Errors  
23 by the Institute of Medicine (42) and carry important legal and clinical implications:

- 24 - Legally, misdiagnoses are the most common source of malpractice claims or litigation.<sup>(43)</sup>
- 25 - Clinically, missed fractures in radiographs have consequences such as malunion with restricted  
26 range of motion, posttraumatic osteoarthritis, and joint collapse (44)

1 Physicians commonly tackle these implications by performing confirmative CT and/or MR studies,  
2 inevitably resulting in a delay in diagnosis, increase in costs and potentially also higher exposure to  
3 radiation. The delays can range from hours to months, resulting in poorer clinical outcomes.  
4 (10,11)

5 A commonly mentioned rebuttal for the implementation of AI based algorithms is the so called  
6 “black box” problem, where the clinician is blinded to the “reasoning” behind the model’s prediction.  
7 (45) Visualization techniques such as heatmaps could improve the acceptance of fracture  
8 detection systems in the clinical practice.

9

## 10 **Conclusion**

11 This study found that our AI model can accurately identify anomalies suggestive of thoracolumbar  
12 vertebral fractures in sagittal radiographs. Specifically, an adapted version based on ResNet18  
13 achieved a similar performance compared to other models, and those reported of expert surgeons  
14 and radiologists. Additionally, it also highlighted a human error made during the annotation  
15 process. Applying this AI model to minimize diagnostic errors in fracture detection in sagittal  
16 radiographs of the TL vertebra seems plausible.

17

## 18 **Conflict of interest**

19 The named authors have no conflict of interest, financial or otherwise.

20

## 21 **References**

- 22 1. Liu B, Zhu Y, Liu S, Chen W, Zhang F, Zhang Y. National incidence of traumatic spinal  
23 fractures in China: Data from China National Fracture Study. *Medicine (Baltimore)*. 2018  
24 Aug;97(35):e12190–e12190.
- 25 2. Hu R, Mustard CA, Burns C. Epidemiology of Incident Spinal Fracture in a Complete

- 1 Population. *Spine (Phila Pa 1976)* [Internet]. 1996 Feb;21(4):492–9. Available from:  
2 <http://journals.lww.com/00007632-199602150-00016>
- 3 3. Niemi-Nikkola V, Saijets N, Ylipoussu H, Kinnunen P, Pesälä J, Mäkelä P, et al. Traumatic  
4 Spinal Injuries in Northern Finland. *Spine (Phila Pa 1976)*. 2018 Jan;43(1):E45–51.
- 5 4. Daly MC, Patel MS, Bhatia NN, Bederman SS. The Influence of Insurance Status on the  
6 Surgical Treatment of Acute Spinal Fractures. *Spine (Phila Pa 1976)*. 2016 Jan;41(1):E37-  
7 45.
- 8 5. Hasler RM, Exadaktylos AK, Bouamra O, Benneker LM, Clancy M, Sieber R, et al.  
9 Epidemiology and predictors of spinal injury in adult major trauma patients: European cohort  
10 study. *Eur Spine J*. 2011;20(12):2174–80.
- 11 6. Baaj AA, Downes K, Vaccaro AR, Uribe JS, Vale FL. Trends in the treatment of lumbar  
12 spine fractures in the United States: A socioeconomics perspective - Clinical article. *J*  
13 *Neurosurg Spine*. 2011 Oct;15(4):367–70.
- 14 7. Delmas PD, Van Langerijt L De, Watts NB, Eastell R, Genant H, Grauer A, et al.  
15 Underdiagnosis of vertebral fractures is a worldwide problem: The IMPACT study. *J Bone*  
16 *Miner Res*. 2005 Apr;20(4):557–63.
- 17 8. Reinhold M, Audigé L, Schnake KJ, Bellabarba C, Dai LY, Oner FC. AO spine injury  
18 classification system: A revision proposal for the thoracic and lumbar spine. *Eur Spine J*.  
19 2013;22(10):2184–201.
- 20 9. Pizones J, Izquierdo E, Alvarez P, Sanchez-Mariscal F, Zuniga L, Chimeno P, et al. Impact  
21 of magnetic resonance imaging on decision making for thoracolumbar traumatic fracture  
22 diagnosis and treatment. *Eur spine J Off Publ Eur Spine Soc Eur Spinal Deform Soc Eur*  
23 *Sect Cerv Spine Res Soc*. 2011 Aug;20 Suppl 3:390–6.
- 24 10. Rhee PM, Bridgeman A, Acosta JA, Kennedy S, Wang DSY, Sarveswaran J. Lumbar  
25 fractures in adult blunt trauma: Axial and single-slice helical abdominal and pelvic computed

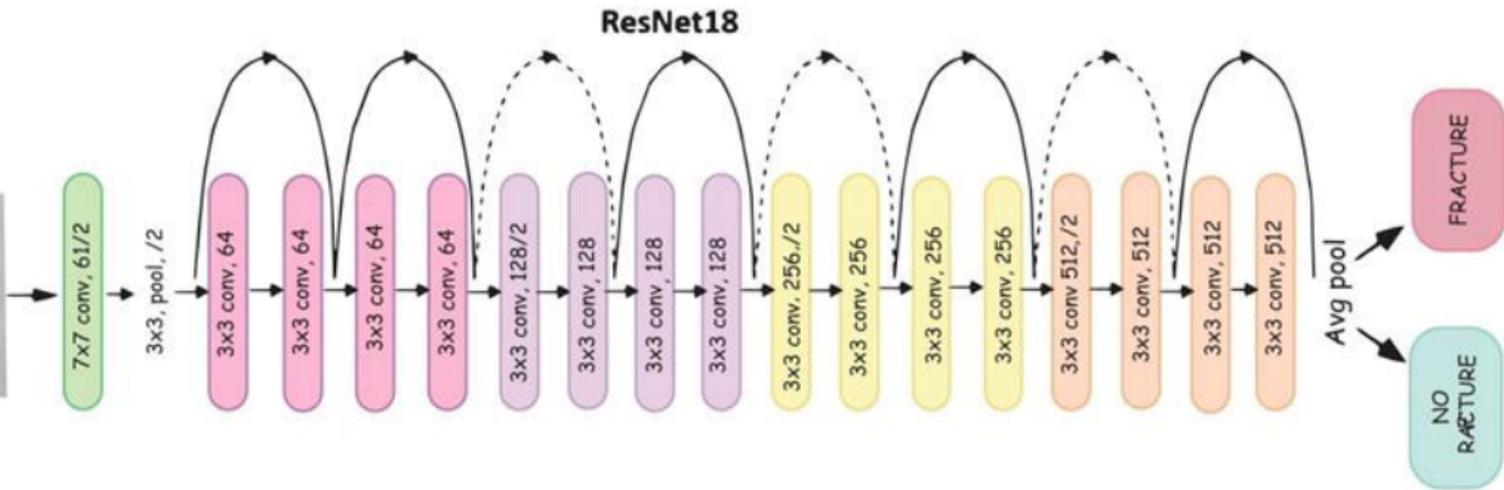
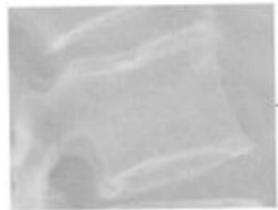
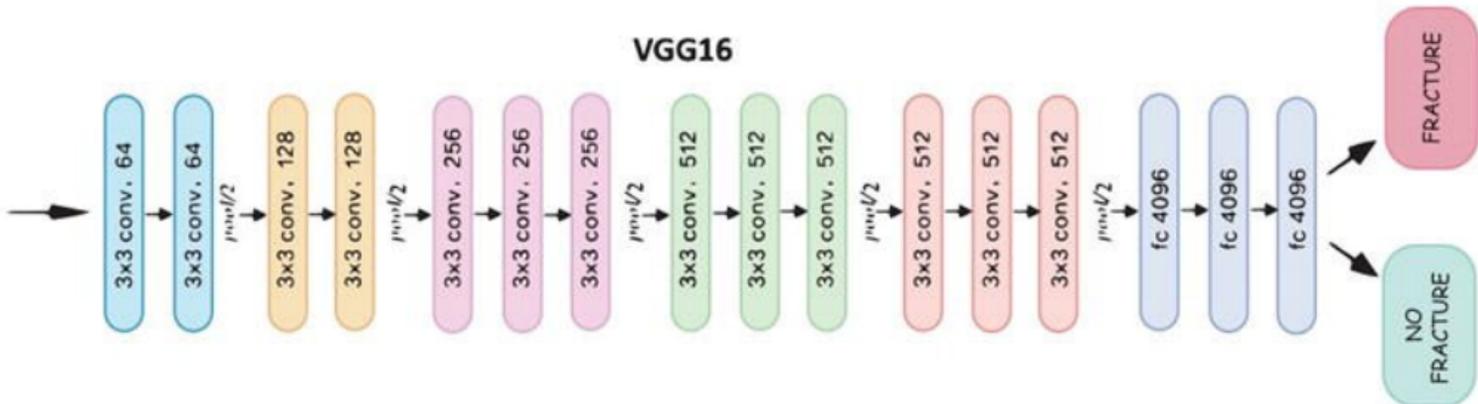
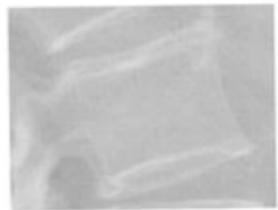
- 1 tomographic scans versus portable plain films. J Trauma [Internet]. 2002 [cited 2021 Mar  
2 19];53(4):663–7. Available from: <https://pubmed.ncbi.nlm.nih.gov/12394863/>
- 3 11. Aso-Escario J, Sebastián C, Aso-Vizán A, Martínez-Quiñones JV, Consolini F, Arregui R.  
4 Delay in diagnosis of thoracolumbar fractures. Orthop Rev (Pavia). 2019 May 23;11(2):47–  
5 52.
- 6 12. Zhang Z, Sejdić E. Radiological images and machine learning: Trends, perspectives, and  
7 prospects [Internet]. Vol. 108, Computers in Biology and Medicine. Elsevier Ltd; 2019 [cited  
8 2020 Apr 24]. p. 354–70. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/31054502>
- 9 13. Savadjiev P, Chong J, Dohan A, Vakalopoulou M, Reinhold C, Paragios N, et al.  
10 Demystification of AI-driven medical image interpretation: past, present and future. Vol. 29,  
11 European Radiology. Springer Verlag; 2019. p. 1616–24.
- 12 14. Hosny A, Parmar C, Quackenbush J, Schwartz LH, Aerts HJWL, Edu HH. Artificial  
13 intelligence in radiology HHS Public Access. Nat Rev Cancer. 2018;18(8):500–10.
- 14 15. Doshi AM, Moore WH, Kim DC, Rosenkrantz AB, Fefferman NR, Ostrow DL, et al.  
15 Informatics Solutions for Driving an Effective and Efficient Radiology Practice.  
16 Radiographics [Internet]. 2018 Oct [cited 2019 Nov 15];38(6):1810–22. Available from:  
17 <http://www.ncbi.nlm.nih.gov/pubmed/30303784>
- 18 16. Ahmed Z, Mohamed K, Zeeshan S, Dong XQ. Artificial intelligence with multi-functional  
19 machine learning platform development for better healthcare and precision medicine.  
20 Database. 2020;2020:1–35.
- 21 17. Zaharchuk G, Gong E, Wintermark M, Rubin D, Langlotz CP. Deep learning in  
22 neuroradiology. Vol. 39, American Journal of Neuroradiology. American Society of  
23 Neuroradiology; 2018. p. 1776–84.
- 24 18. Krizhevsky A, Sutskever I, Hinton GE. ImageNet Classification with Deep Convolutional  
25 Neural Networks [Internet]. [cited 2019 Nov 14]. Available from:

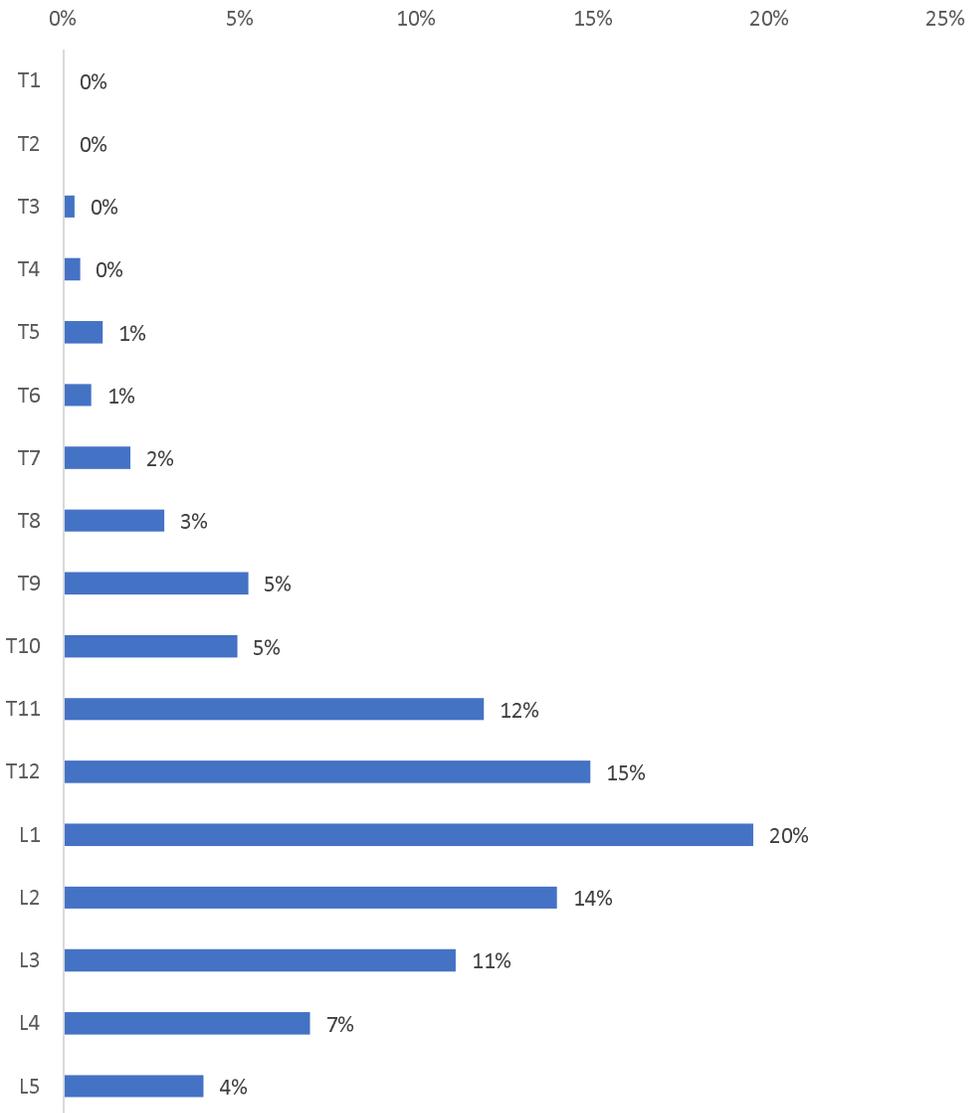
- 1 <http://code.google.com/p/cuda-convnet/>
- 2 19. Gyftopoulos S, Lin D, Knoll F, Doshi AM, Rodrigues TC, Recht MP. Artificial Intelligence in  
3 Musculoskeletal Imaging: Current Status and Future Directions. *AJR Am J Roentgenol*  
4 [Internet]. 2019 Sep [cited 2019 Nov 15];213(3):506–13. Available from:  
5 <https://www.ajronline.org/doi/10.2214/AJR.19.21117>
- 6 20. Jamaludin A, Lootus M, Kadir T, Zisserman A, Urban J, Battié MC, et al. ISSLS PRIZE IN  
7 BIOENGINEERING SCIENCE 2017: Automation of reading of radiological features from  
8 magnetic resonance images (MRIs) of the lumbar spine without human intervention is  
9 comparable with an expert radiologist. *Eur spine J Off Publ Eur Spine Soc Eur Spinal*  
10 *Deform Soc Eur Sect Cerv Spine Res Soc*. 2017 May;26(5):1374–83.
- 11 21. Weng C-H, Wang C-L, Huang Y-J, Yeh Y-C, Fu C-J, Yeh C-Y, et al. Artificial Intelligence for  
12 Automatic Measurement of Sagittal Vertical Axis Using ResUNet Framework. *J Clin Med*.  
13 2019 Nov;8(11).
- 14 22. Vergari C, Skalli W, Gajny L. A convolutional neural network to detect scoliosis treatment in  
15 radiographs. *Int J Comput Assist Radiol Surg*. 2020 Jun;15(6):1069–74.
- 16 23. Maki S, Furuya T, Horikoshi T, Yokota H, Mori Y, Ota J, et al. A Deep Convolutional Neural  
17 Network With Performance Comparable to Radiologists for Differentiating Between Spinal  
18 Schwannoma and Meningioma. *Spine (Phila Pa 1976)*. 2020 May;45(10):694–700.
- 19 24. Chmelik J, Jakubicek R, Walek P, Jan J, Ourednicek P, Lambert L, et al. Deep convolutional  
20 neural network-based segmentation and classification of difficult to define metastatic spinal  
21 lesions in 3D CT data. *Med Image Anal*. 2018 Oct;49:76–88.
- 22 25. Galbusera F, Casaroli G, Bassani T. Artificial intelligence and machine learning in spine  
23 research. *JOR SPINE [Internet]*. 2019 Mar 5;2(1):e1044. Available from:  
24 <https://onlinelibrary.wiley.com/doi/abs/10.1002/jsp2.1044>
- 25 26. Murata K, Endo K, Aihara T, Suzuki H, Sawaji Y, Matsuoka Y, et al. Artificial intelligence for

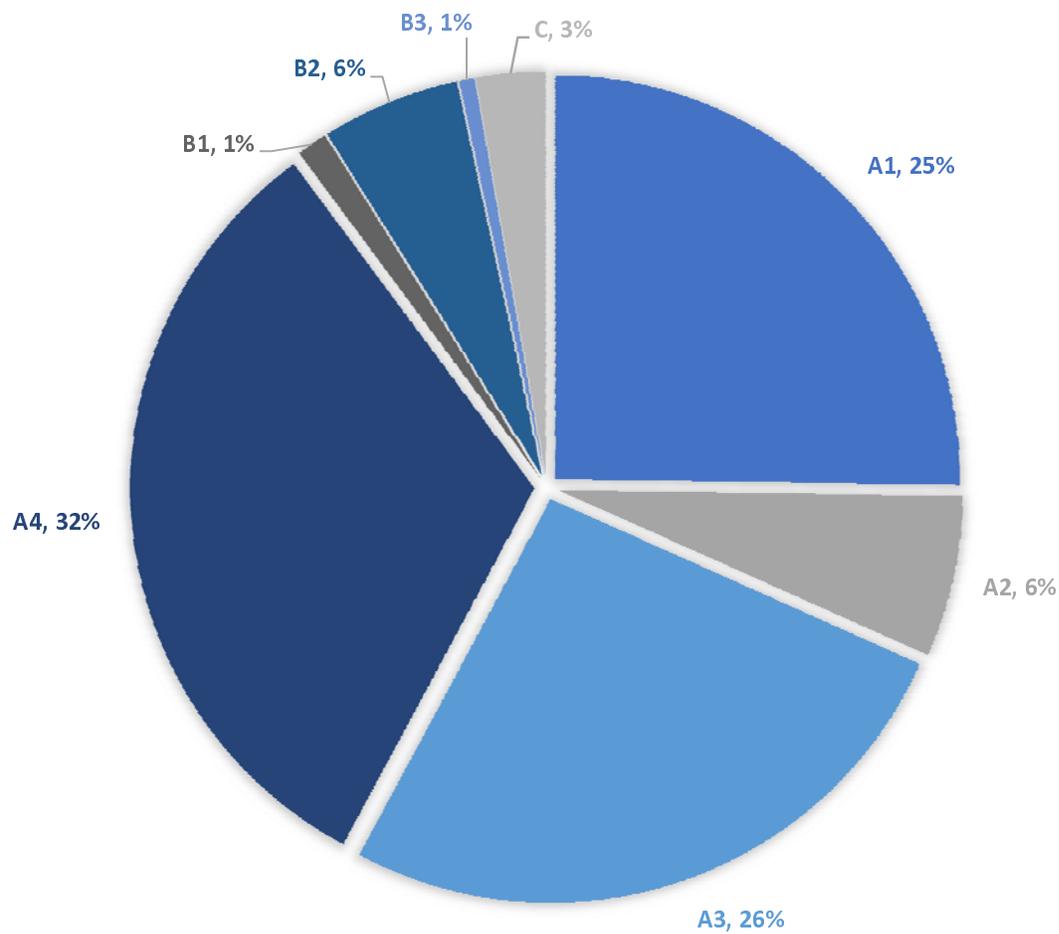
- 1 the detection of vertebral fractures on plain spinal radiography. *Sci Rep* [Internet].  
2 2020;10(1):1–8. Available from: <https://doi.org/10.1038/s41598-020-76866-w>
- 3 27. Simonyan K, Zisserman A. Very deep convolutional networks for large-scale image  
4 recognition. *3rd Int Conf Learn Represent ICLR 2015 - Conf Track Proc.* 2015;1–14.
- 5 28. He K, Zhang X, Ren S, Sun J. Deep residual learning for image recognition. *Proc IEEE  
6 Comput Soc Conf Comput Vis Pattern Recognit.* 2016;2016-Decem:770–8.
- 7 29. Deng J, Dong W, Socher R, Li L-J, Kai Li, Li Fei-Fei. ImageNet: A large-scale hierarchical  
8 image database. 2010;(May 2014):248–55.
- 9 30. Shin HC, Roth HR, Gao M, Lu L, Xu Z, Nogues I, et al. Deep Convolutional Neural Networks  
10 for Computer-Aided Detection: CNN Architectures, Dataset Characteristics and Transfer  
11 Learning. *IEEE Trans Med Imaging.* 2016;35(5):1285–98.
- 12 31. Panigrahi S, Nanda A, Swarnkar T. A Survey on Transfer Learning. *Smart Innov Syst  
13 Technol.* 2021;194:781–9.
- 14 32. Sidey-Gibbons JAM, Sidey-Gibbons CJ. Machine learning in medicine: a practical  
15 introduction. *BMC Med Res Methodol* [Internet]. 2019 Mar 19 [cited 2020 Aug 13];19(1):64.  
16 Available from: [https://bmcmmedresmethodol.biomedcentral.com/articles/10.1186/s12874-  
17 019-0681-4](https://bmcmmedresmethodol.biomedcentral.com/articles/10.1186/s12874-019-0681-4)
- 18 33. Olczak J, Fahlberg N, Maki A, Razavian AS, Jilert A, Stark A, et al. Artificial intelligence for  
19 analyzing orthopedic trauma radiographs: Deep learning algorithms—are they on par with  
20 humans for diagnosing fractures? *Acta Orthop* [Internet]. 2017 Nov 2 [cited 2020 Aug  
21 13];88(6):581–6. Available from: [/pmc/articles/PMC5694800/?report=abstract](https://pubmed.ncbi.nlm.nih.gov/30000000/)
- 22 34. Adams M, Chen W, Holcdorf D, McCusker MW, Howe PDL, Gaillard F. Computer vs human:  
23 Deep learning versus perceptual training for the detection of neck of femur fractures. *J Med  
24 Imaging Radiat Oncol.* 2019;63(1):27–32.
- 25 35. Chung SW, Han SS, Lee JW, Oh KS, Kim NR, Yoon JP, et al. Automated detection and

- 1 classification of the proximal humerus fracture by using deep learning algorithm. *Acta*  
2 *Orthop*. 2018;89(4):468–73.
- 3 36. Kim DH, MacKinnon T. Artificial intelligence in fracture detection: transfer learning from deep  
4 convolutional neural networks. *Clin Radiol* [Internet]. 2018;73(5):439–45. Available from:  
5 <https://doi.org/10.1016/j.crad.2017.11.015>
- 6 37. Urakawa T, Tanaka Y, Goto S, Matsuzawa H, Watanabe K, Endo N. Detecting  
7 intertrochanteric hip fractures with orthopedist-level accuracy using a deep convolutional  
8 neural network. *Skeletal Radiol*. 2019;48(2):239–44.
- 9 38. Yang S, Yin B, Cao W, Feng C, Fan G, He S. Diagnostic accuracy of deep learning in  
10 orthopaedic fractures: a systematic review and meta-analysis. *Clin Radiol*.  
11 2020;75(9):713.e17-713.e28.
- 12 39. Bluemke DA, Moy L, Bredella MA, Ertl-Wagner BB, Fowler KJ, Goh VJ, et al. Assessing  
13 radiology research on artificial intelligence: A brief guide for authors, reviewers, and readers-  
14 from the Radiology Editorial Board. Vol. 294, *Radiology*. Radiological Society of North  
15 America Inc.; 2020. p. 487–9.
- 16 40. Chung SW, Han SS, Lee JW, Oh KS, Kim NR, Yoon JP, et al. Automated detection and  
17 classification of the proximal humerus fracture by using deep learning algorithm. *Acta*  
18 *Orthop* [Internet]. 2018 Jul 4 [cited 2020 Aug 13];89(4):468–73. Available from:  
19 [/pmc/articles/PMC6066766/?report=abstract](https://pubmed.ncbi.nlm.nih.gov/3207237/)
- 20 41. Espinosa JA, Nolan TW. Reducing errors made by emergency physicians in interpreting  
21 radiographs: Longitudinal study. *Br Med J* [Internet]. 2000 Mar 18 [cited 2021 Mar  
22 20];320(7237):737–40. Available from: <https://pubmed.ncbi.nlm.nih.gov/10720354/>
- 23 42. Balogh EP, Miller BT, Ball JR. Improving diagnosis in health care [Internet]. *Improving*  
24 *Diagnosis in Health Care*. National Academies Press; 2016 [cited 2021 Mar 22]. 1–472 p.  
25 Available from: <https://www.ncbi.nlm.nih.gov/books/NBK338596/>

- 1 43. Kachalia A, Gandhi TK, Puopolo AL, Yoon C, Thomas EJ, Griffey R, et al. Missed and  
2 Delayed Diagnoses in the Emergency Department: A Study of Closed Malpractice Claims  
3 From 4 Liability Insurers. *Ann Emerg Med* [Internet]. 2007 Feb [cited 2021 Mar  
4 20];49(2):196–205. Available from: <https://pubmed.ncbi.nlm.nih.gov/16997424/>
- 5 44. Lindsey R, Daluiski A, Chopra S, Lachapelle A, Mozer M, Sicular S, et al. Deep neural  
6 network improves fracture detection by clinicians. *Proc Natl Acad Sci U S A* [Internet]. 2018  
7 Nov 6 [cited 2020 Aug 13];115(45):11591–6. Available from:  
8 <https://pubmed.ncbi.nlm.nih.gov/30348771/>
- 9 45. Fazal MI, Patel ME, Tye J, Gupta Y. The past, present and future role of artificial intelligence  
10 in imaging. *Eur J Radiol* [Internet]. 2018;105:246–50. Available from:  
11 <https://doi.org/10.1016/j.ejrad.2018.06.020>
- 12







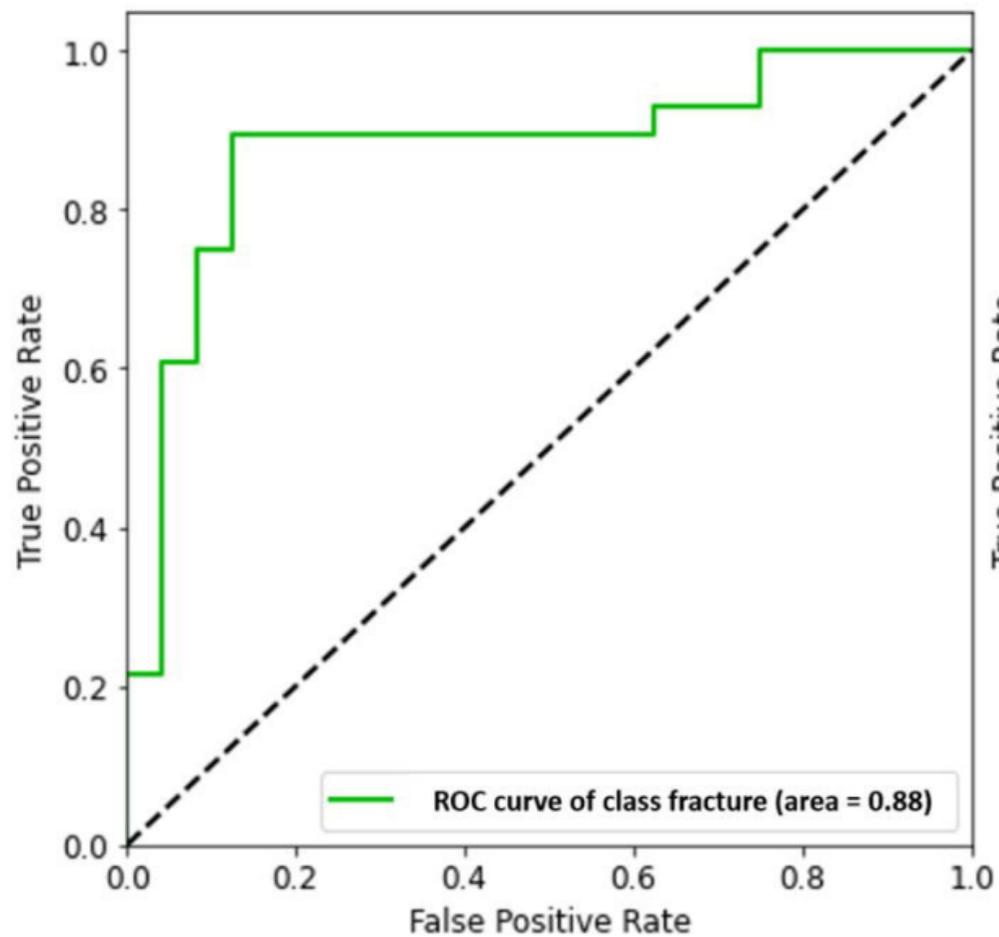
ResNet18

Ground Truth	no fracture	TN 25 48.1%	FP 3 5.8%
	fracture	FN 3 5.8%	TP 21 40.4%
		no fracture	fracture
		Model	

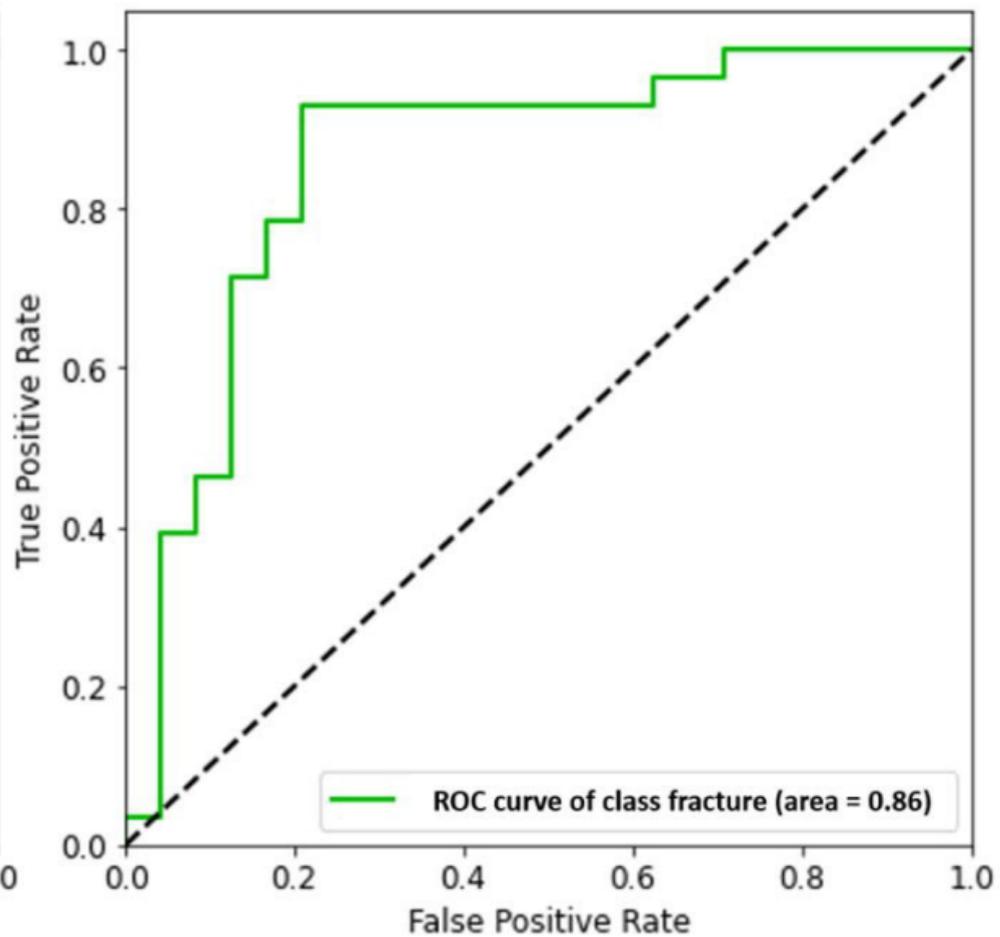
VGG16

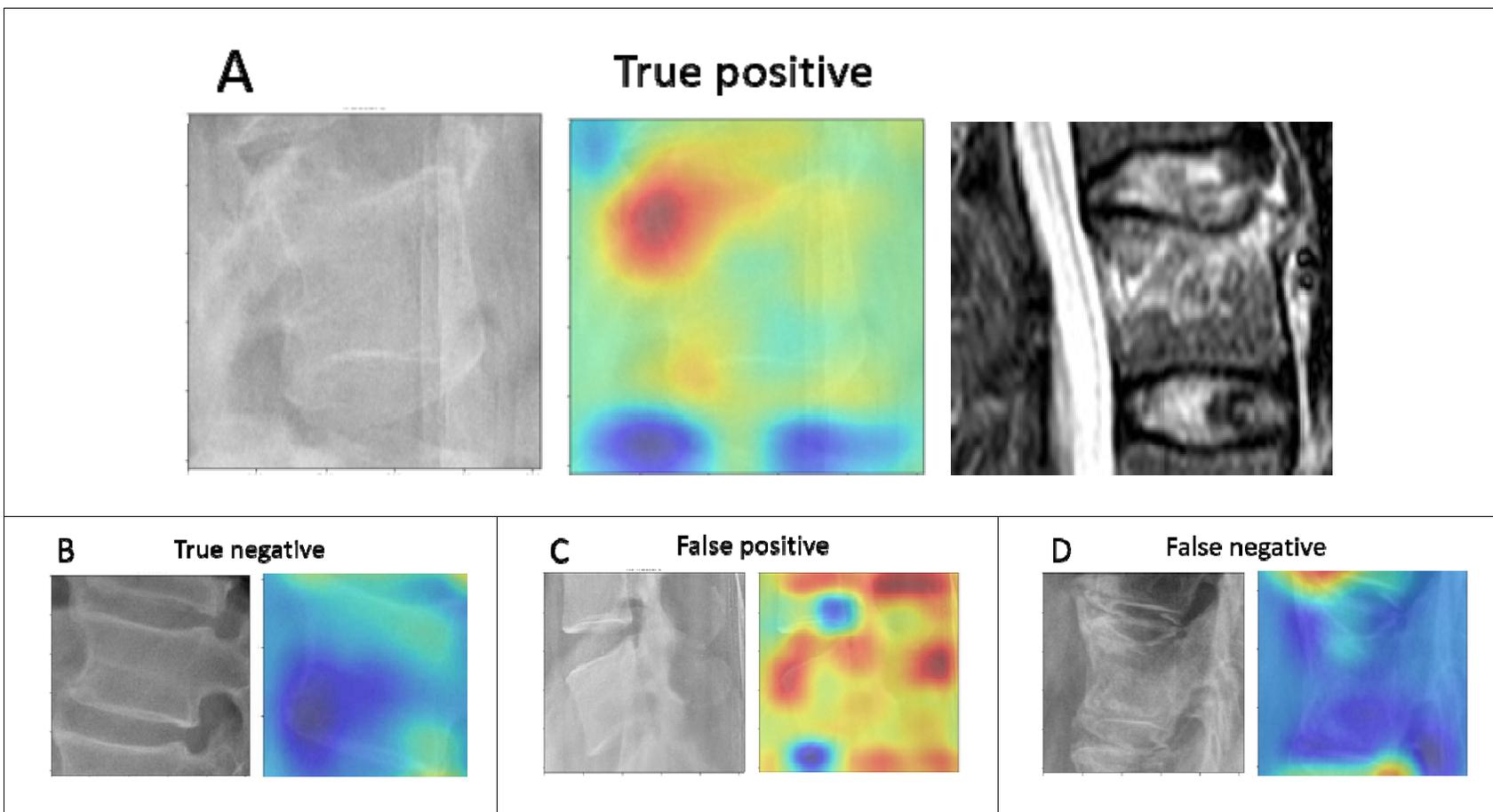
Ground Truth	no fracture	TN 25 48.1%	FP 5 9.6%
	fracture	FN 3 5.8%	TP 19 36.5%
		no fracture	fracture
		Model	

ROC ResNet18



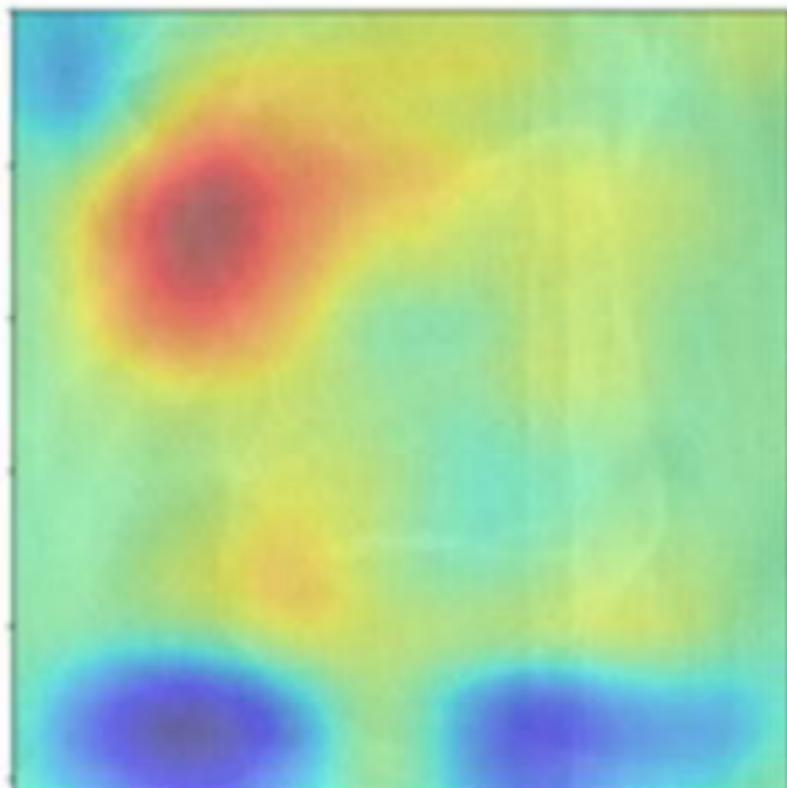
ROC VGG16





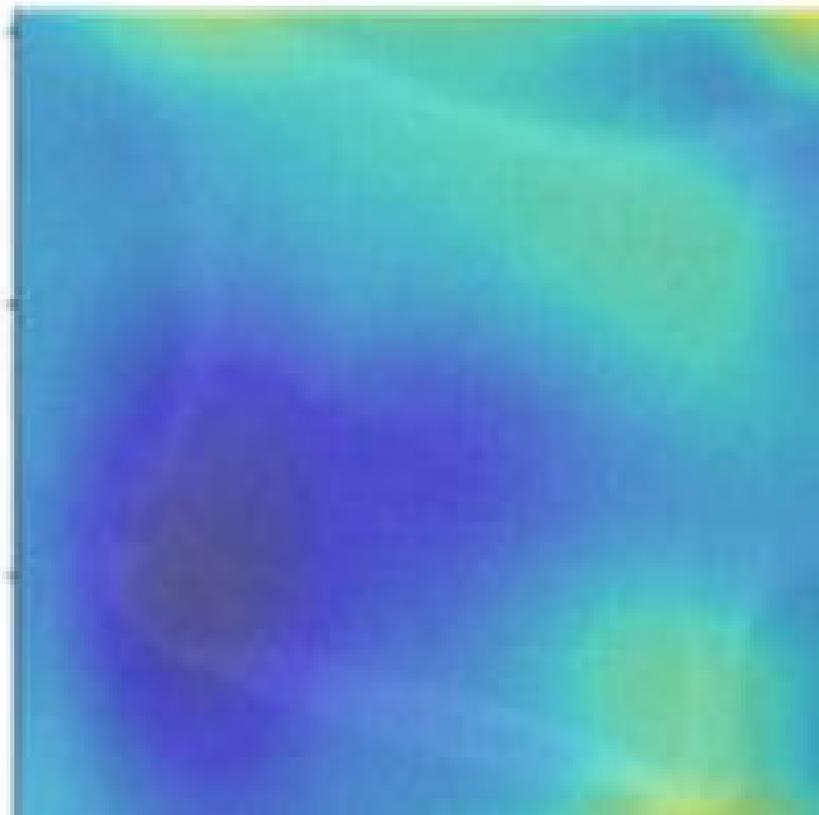
A

True positive



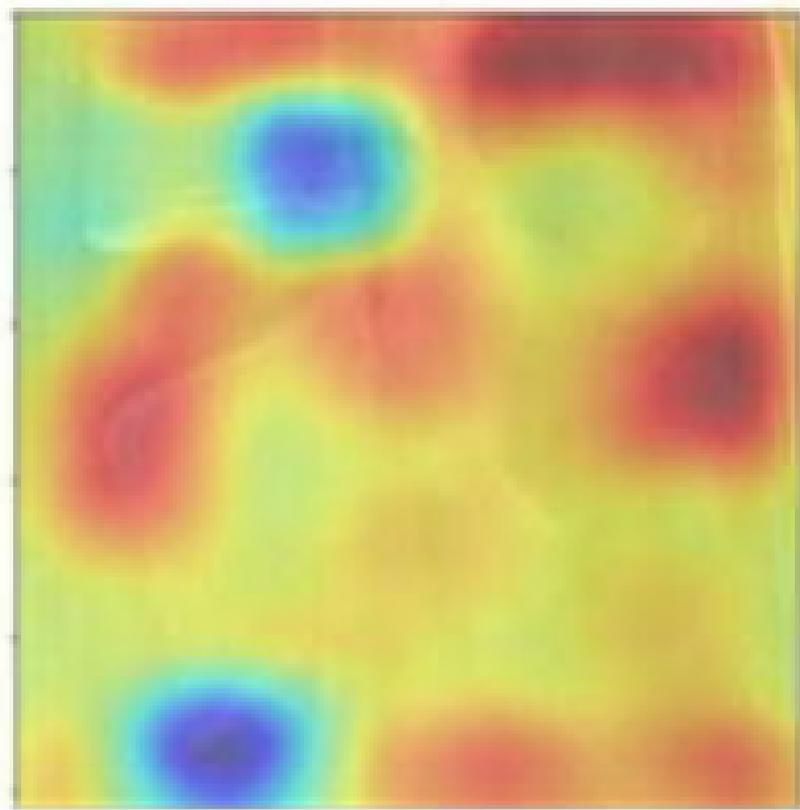
**B**

True negative



C

False positive



D

False negative

