

Diabetic Foot Ulcer Grand Challenge 2021: Evaluation and Summary

Bill Cassidy¹[0000-0003-3741-8120], Connah Kendrick¹[0000-0002-3623-6598], Neil D. Reeves²[0000-0001-9213-4580], Joseph M. Pappachan³[0000-0003-0886-5255], Claire O'Shea⁴, David G. Armstrong⁵[0000-0003-1887-9175], and Moi Hoon Yap¹[0000-0001-7681-4287]

- ¹ Department of Computing and Mathematics, Manchester Metropolitan University, Manchester M1 5GD, UK
² Musculoskeletal Science and Sports Medicine, Manchester Metropolitan University, Manchester M1 5GD, UK
³ Lancashire Teaching Hospitals NHS Foundation Trust, Preston, PR2 9HT, UK
⁴ Waikato District Health Board, Hamilton 3240, New Zealand
⁵ Southwestern Academic Limb Salvage Alliance (SALSA), Department of Surgery, Keck School of Medicine of University of Southern California, Los Angeles, California, USA
B.Cassidy@mmu.ac.uk, M.Yap@mmu.ac.uk

Abstract. Diabetic foot ulcer classification systems use the presence of wound infection (bacteria present within the wound) and ischaemia (restricted blood supply) as vital clinical indicators for treatment and prediction of wound healing. Studies investigating the use of automated computerised methods of classifying infection and ischaemia within diabetic foot wounds are limited due to a paucity of publicly available datasets and severe data imbalance in those few that exist. The Diabetic Foot Ulcer Challenge 2021 provided participants with a more substantial dataset comprising a total of 15,683 diabetic foot ulcer patches, with 5,955 used for training, 5,734 used for testing and an additional 3,994 unlabelled patches to promote the development of semi-supervised and weakly-supervised deep learning techniques. This paper provides an evaluation of the methods used in the Diabetic Foot Ulcer Challenge 2021, and summarises the results obtained from each network. The best performing network was an ensemble of the results of the top 3 models, with a macro-average F1-score of 0.6307.

1 Introduction

Diabetic foot ulcers (DFU) are one of the most serious complications that can result from diabetes, and often lead to amputation of all or part of a limb if not met with timely treatment [1,2]. Early detection of DFU, together with accurate screening for infection and ischaemia can help in early treatment and avoidance of further serious complications including amputation. In previous studies, various researchers [3,4,5,6,7] have achieved high accuracy in automated detection of DFUs with machine learning algorithms. A number of widely used clinical DFU

classification systems are currently in use, such as Wagner [8], University of Texas [9,10], and SINBAD Classification [11], which include information on the site of the DFU, area, depth, presence of neuropathy, ischaemia and infection. We focus on ischaemia and infection, which are key features of DFU classification systems and important clinical determinants for effective treatment and healing. This focus is consistent with the evolution of threatened limb classification systems, such as the Wound, Ischemia, and foot Infection (WIFI) classification which is used to predict the risk of amputation in patients diagnosed with critical limb ischemia [12,13].

Recognition of infection and ischaemia are key determinate factors that predict the healing progress of DFU and risk of amputation. Ischaemia develops due to lack of arterial inflow to the foot, that results in spontaneous necrosis of the most poorly perfused tissues (gangrene), which may ultimately require amputation of part of the foot or leg. In previous studies, it is estimated that patients with critical limb ischaemia have a three-year limb loss rate of approximately 40% [14]. Patients with an active DFU, particularly those with ischaemia or gangrene, should also be examined for the presence of infection. Approximately, 56% of DFU become infected and 20% of DFU infections lead to amputation of foot or limb [15,16,17]. In one recent study, 785 million patients with diabetes in the US between 2007 and 2013 suggested that DFU and associated infections constitute a powerful risk factor for emergency department visits and hospital admission [18]. Due to high risks of infection and ischaemia associated with DFU amputation [12], timely and accurate recognition of infection and ischaemia in DFU with cost-effective machine learning methods is an important step towards the development of a complete computerised DFU assessment system.

In current practice, DFU assessment is conducted in foot clinics and hospitals by podiatrists and diabetes physicians. To determine appropriate treatment, a vascular assessment is performed for ischaemia and the wound is assessed for clinical evidence of infection and wound tissue sent for microbiological culture. Van Netten et al. [19] found that clinicians achieved low validity and reliability for remote assessment of DFU in foot images. Hence, it is clear that analysing these conditions from images is extremely difficult even by experienced podiatrists. Patient experiences may be different, however. Swerdlow et al. [20], instituted a “foot selfie” programme and found overall high levels of patient engagement. Limited research exists using computerised methods to automate the monitoring of DFU using foot photographs [21]. This is due to the lack of availability of datasets with clinical labelling for research purposes [22].

Motivated by technological advancements in medical imaging [23,24,25,26,27], where machine learning algorithms performed better than experienced clinicians, Goyal et al. [28] analysed the performance of machine learning algorithms on the recognition of ischaemia and infection on DFUs. They proved that deep learning methods outperformed conventional machine learning methods on a small dataset (1,459 images) and proposed an ensemble Convolutional Neural Network (CNN) approach for ischaemia and infection recognition. Although they achieved high accuracy in ischaemia recognition, there were a number of limitations to

their method: 1) the proposed binary classification ensemble CNN method detected one class at a time, which was not capable of detecting co-occurrence of infection and ischaemia; 2) the dataset was small and cannot be generalised; 3) the dataset was highly imbalanced, with infection cases significantly outnumbering ischaemic cases; and 4) the recognition rate of infection was 73%, which requires substantial work to improve accuracy. To address these issues, Yap et al. [29] introduced the DFUC2021 datasets, which consist of 4,474 clinically annotated images, together with DFU patches with the label of infection, ischaemia, both infection and ischaemia and none of those conditions (control). Since the release of the DFUC2021 datasets on the 15th April 2021, they have been shared with 51 institutions from 25 countries. Figure 1 illustrates the distribution of researchers using the DFUC2021 datasets by country, showing that the majority of users originate from the United States, China, India and Brazil.

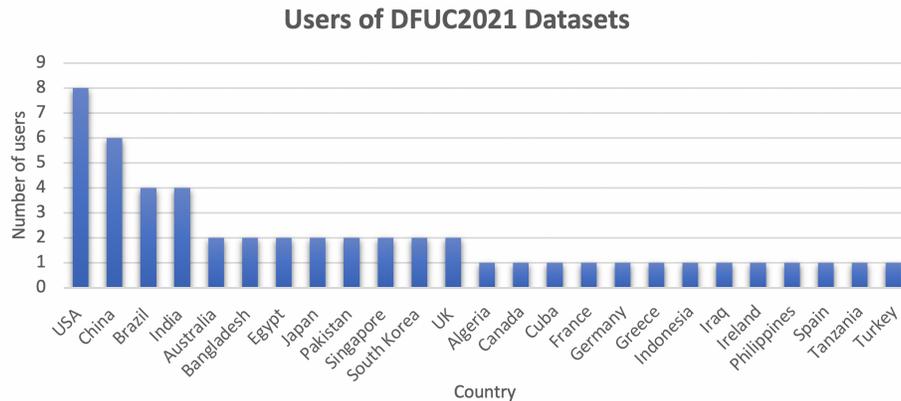


Fig. 1. Distribution of researchers using the DFUC2021 datasets and their country of origin.

2 Methodology

This section summarises the datasets used for DFUC2021, the performance metrics and analysis of the methods proposed by the participants for challenge.

2.1 Datasets and Ground Truth

The previous publicly available dataset created by Goyal et al. [28] consists of 1,459 DFUs: 645 with infection, 24 with ischaemia, 186 with infection and ischaemia, and 604 control DFU (presence of DFU, but without infection or ischaemia). The DFUC2021 dataset is the largest publicly available dataset with

DFU pathology labels which consists of 1,703 ulcers with infection, 152 ulcers with ischaemia, 372 ulcers with both conditions, 1703 control DFU and an additional 1,337 unlabelled DFU. The ground truth was produced by two healthcare professionals who specialise in diabetic wounds and ulcers. The instruction for annotation was to label each ulcer with ischaemia, and/or infection, or none. The patient medical record was used to validate the labels. To increase the number of DFU patches for deep learning algorithms, we used natural augmentation [28] and generated a total of 15,683 DFU patches, which consists of 11,689 labelled patches and 3,994 unlabelled patches. For the labelled DFU patches, the train, validation and test split is: 4,799 patches for the training set, 1,156 patches for the validation set, and 5,734 patches for the testing set. The detailed split for each pathology is presented in Table 1. As shown in Table 1, the number of patches for ischaemia, infection and ischaemia are relatively low when compared to the other classes.

Table 1. DFUC2021 dataset distribution for training (4,799 patches) and validation (1,156 patches) after natural augmentation.

	Infection	Ischaemia	Infection and Ischaemia	None	Total
Train	2074	179	483	2063	4799
Validation	481	48	138	489	1156

Figure 2 illustrates DFU patches of four conditions. As these ulcers exhibit variability within a single condition and similarity between different conditions, the DFUC2021 dataset presents a significant challenge for computer vision and machine learning methods in the recognition of infection and ischaemia.

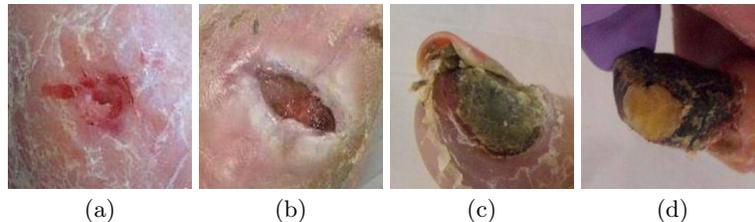


Fig. 2. Illustration of DFU patches from the DFUC2021 dataset. From left to right: (a) control DFU, (b) DFU with infection, (c) DFU with ischaemia and (d) DFU with both infection and ischaemia.

2.2 Performance Metrics

We compared the performance of the deep learning networks on recognition of infection and ischaemia using precision, recall, F1-score and area under the

Receiver Operating Characteristics Curve (AUC). For the performance in multi-label classification, due to class imbalance inherent within the DFUC2021 dataset, the performance will be reported in macro-average. Macro-average is used in imbalanced multi-class settings as it provides equal emphasis on minority classes [30]. To compute macro-average F1-score, first we obtain all the True Positives (TP), False Positives (FP) and False Negatives (FN) for each class i (of n classes), and their respective F1-scores. The Macro-F1 is determined by averaging the per-class F1-scores, $F1_i$:

$$F1_i = \frac{2 \cdot TP_i}{2 \cdot TP_i + FP_i + FN_i} \quad (1)$$

$$\text{Macro-F1} = \frac{1}{N} \sum_i F1_i \quad (2)$$

where $i = 1, \dots, N$ represents the i -th class and N is the total number of classes, in this case, $N = 4$. For completeness, we also compared the performance of the algorithms with macro-average AUC.

2.3 Analysis of the Proposed Methods

In this section, we detail the methods used by the top 10 entries for DFUC2021.

The method ranked 10th in the challenge, submitted by Ye Hai, proposes two classifiers. The first classifier is used to detect control cases, while the second classifier is used to detect the other three categories - infection, ischaemia and both infection and ischaemia. The group tested a variety of classification networks, such as ResNet, ViT, DenseNet and SENet, and found that SENet34 provided the best results for both classifiers.

The method ranked 9th in the challenge, submitted by Weilun Wang, proposes a texture classification model which used SE-DenseNet (Squeeze and Excitation Densely Connected Convolutional Network). SE-DenseNet combines the advantages of DenseNet and SENet, which uses multi-dimensional feature information, strengthening the transmission of deep information and enhancing the learning and expression ability of the deep network through a "feature recalibration" strategy. Further, the network is also able to slow down the attenuation of errors in each hidden layer, ensuring the stability of gradient weight information and avoiding the disappearance of gradient through the reverse conduction mechanism of the network itself, improving network performance [31]. This approach does not require a very deep model, so networks such as DenseNet121 and EfficientNetB0, which contain over 100 convolution layers, were not used. To determine the optimum model depth required for this scenario, Wang performed several experiments. First, the 5,955 training samples were split into 10 subsets, followed by 10-fold cross validation. Next, 9-fold samples were randomly up-sampled and used as training samples in each sub-experiment, with the remaining 1-fold samples used for testing.

The method ranked 8th in DFUC2021, submitted by Das et al., proposed a prediction level ensembling. This submission utilised DenseNet121 and EfficientNetB0 models pretrained using ImageNet. The convolutional layers are taken as

proposed in the original work, however, the fully connected layers are set as FC(4096), FC(4096), FC(1000) and FC(4), which all use ReLU activations except the final softmax based prediction layer. The configuration of FC layers is motivated by the original VGG16 architecture [32]. The Softmax predictions from both networks are averaged to obtain a prediction level ensembling, providing a final prediction. Figure 3 shows an overview of the network configuration⁶.

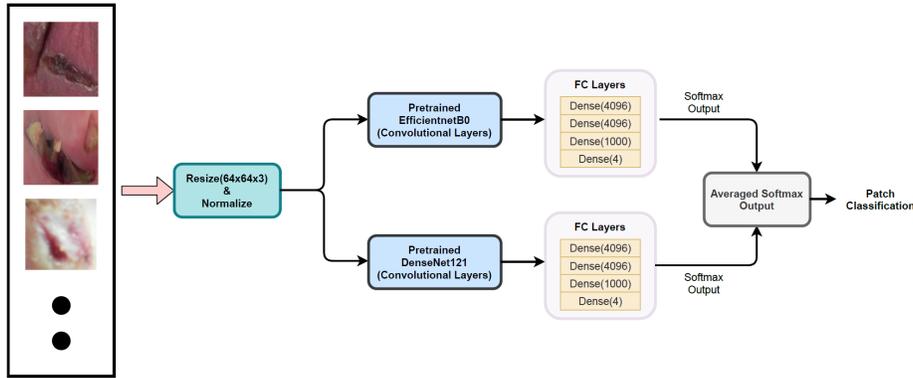


Fig. 3. Illustration of the prediction level ensembling approach used by Das et al., which achieved 8th place in DFUC2021.

The method ranked 7th for DFUC2021, submitted by Chuantao Xie, utilises EfficientNetB0 and DenseNet121, both pretrained using ImageNet. EfficientNetB0 was chosen for its additive feature fusion, while DenseNet121 was chosen for its concatenated feature fusion. The proposed method replaces the layers after the main structure of the CNN, leaving the rest of the network unchanged, including the network structure that proceeds the global average pooling layer which is connected using a parallel structure. One of the branches of the parallel network structure includes convolution, batch normalization, activation function, global average pooling and full connection, followed by the class prediction. Finally, the predicted results of the parallel network structure are concatenated, providing the full connection prediction.

The method which placed 6th for DFUC2021, submitted by Chen et al., utilises an ensemble approach using DenseNet121 and EfficientNet pretrained on ImageNet with a frozen output layer connected to a global average pooling layer. Concatenated integration was implemented with two inputs and one output. The fully connected output layer of the pretrained network was replaced by

⁶ reproduced with permission from Sujit Kumar Das, Department of CSE, National Institute of Technology, Silchar 788010, Assam, India

a new four-class SoftMax layer. Figure 4 shows an overview of the network configuration⁷.

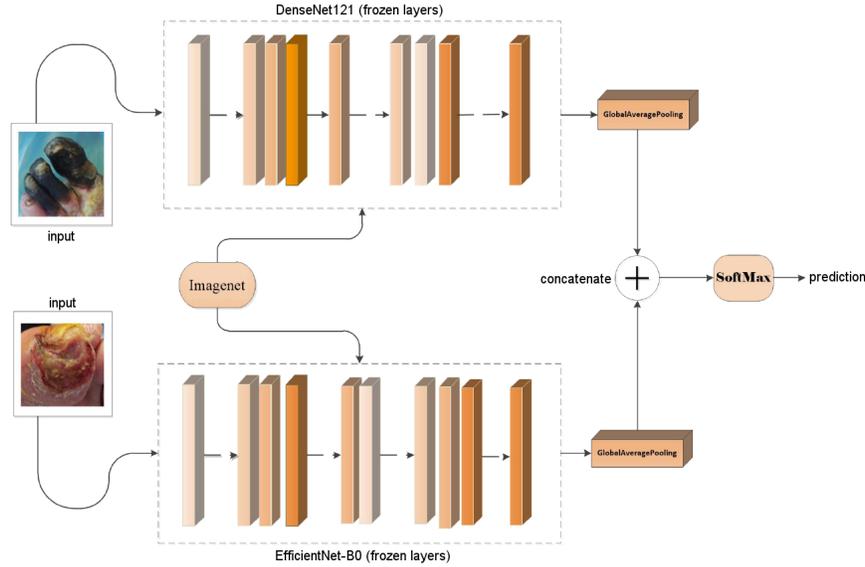


Fig. 4. Illustration of the ensemble CNN approach used by Chen et al., which achieved 6th place in DFUC2021.

The method ranked at 5th place in the challenge, submitted by Güley et al., leveraged the GenerAlly Nuanced Deep Learning Framework (GaNDLF) to achieve multi-class DFU wound classification. GaNDLF enables various machine learning (ML) and artificial intelligent (AI) workloads, including segmentation, regression, classification, and synthesis. This is achieved using a range of imaging modalities, such as RGB, radiographic and histopathologic imaging techniques. Three VGG architectures were trained (VGG11, VGG16 and VGG19) using the DFUC2021 dataset. VGG was selected due to its use of very small convolutions which utilise spacial padding to preserve features from the input image. A total of 5 max-pooling operations were used over a $2 \times N^2$ window size, with a stride of 2 to ensure that image dimensions were halved after each max-pooling operation. ReLU activation with global average pooling and two drop-out layers with a penultimate linear layer were used for the classifier. A Softmax layer forms the last layer to provide a final classification.

The method submitted by Qayyum et al, which ranked 4th in the challenge, utilised Vision Transformers (ViTs) to perform DFU classification. ViTs inherently reduce inductive biases, such as translation variances and locality, which

⁷ reproduced with permission from Donghui Lv, Yuqian Chen, School of Communication and Information Engineering, Shanghai University, Shanghai 200444, China

are present in most other CNN architectures. This solution used pretrained ViTs and fine-tuned them on the DFUC2021 dataset. The features obtained from the last network layer from each ViT were concatenated pairwise, followed by a fully connected layer to concatenate features from individual ViTs before being passed to the final classifier. To address the issue of imbalanced class distribution within the DFUC2021 dataset, average weighted sampling was used, which was shown to improve experimental results.

The method submitted by Ahmed et al., which ranked 3rd in DFUC2021, fine-tuned EfficientNet B0-B6, Resnet-50, and Resnet-101 (pretrained on ImageNet) on the DFUC2021 training set and proposed a custom activation layer using Bias Adjustable Softmax. This Softmax-based activation layer is used to handle class imbalance inherent within the dataset. Their initial experiments used weighted categorical cross entropy but found no significant impact on performance. They found that the use of class weights in the loss function resulted in trained networks showing a bias towards control and infection cases. To address this problem, a novel method was introduced to adjust the skew of probabilities for each class to adjust the bias at inference level.

The method which ranked 2nd place in DFUC2021, submitted by Bloch et al., utilised an ensemble of EfficientNets with a semi-supervised training strategy involving pseudo-labeling for unlabeled images. Their main contribution was the use of Conditional GANs (pix2pixHD) to generate synthetic DFU images to address class imbalance. To achieve this, they created edge masks to indicate regions of interest on the DFUC2021 dataset images.

The winning entry for DFUC2021, submitted by Galdran et al., compared established CNNs (ResNeXt50 and EfficientNet-B3 pretrained on ImageNet) with a Vision Transformer (ViT) and Data-efficient Image Transformers (DeiT) for DFU multi-classification. They also demonstrated how the Sharpness-Aware Minimization (SAM) optimisation algorithm significantly improves the generalisation capability of both traditional CNNs and ViTs in this domain compared to standard Stochastic Gradient Descent (SGD). SAM seeks parameters that lie in neighbourhoods that have uniformly low loss, which results in a min-max optimisation problem on which gradient descent can be performed efficiently, as shown in Algorithm 1. This method was developed to address the problem of heavily overparameterised models where training loss values do not always reflect how well the model generalises. SAM has also been shown to improve robustness against label noise [33]. Their winning entry utilised a linear combination of predictions extracted from BiT-ResNeXt50 (derived from Big Image Transfer) and EfficientNet-B3 models trained on different data folds. This winning submission achieved the highest F1-score (62.16), AUC (88.55) and recall (65.22) measures for DFUC2021.

Algorithm 1 Pseudocode for the SAM algorithm used by the winning entry for DFUC2021, originally proposed by Foret et al. [33].

Input Training set, loss function, batch size, step size, neighborhood size.

Output Model trained with SAM.

```
1: Initialise weights  $w_0$ ,  $t = 0$ 
2: while not converged do
3:   Sample batch  $\beta = \{(x_1, y_1), \dots, (x_b, y_b)\}$ 
4:   Compute gradient  $\nabla_w L_\beta(w)$  of the batch's training loss
5:   Compute  $\hat{\epsilon}(w)$ 
6:   Compute gradient approximation for the SAM objective:  $g = \nabla_w L_\beta(w)|_{w+\hat{\epsilon}(w)}$ 
7:   Update weights:  $w_{t+1} = w_t - \eta g$ 
8:    $t = t + 1$ 
9: return  $w_t$ 
```

3 Results and Discussion

For DFUC2021, there were 500 submissions for the validation stage and 28 submissions for the testing stage. Table 2 shows a summary of the top 10 best performing networks submitted to the DFUC2021 test leaderboard.

3.1 Analysis on the Top-3 Results

In this section, we conduct a statistical analysis of the top-3 results from DFUC2021. Galdran et al. achieved the best F1-score for detection of control cases (0.76), which was an improvement of 0.03 on the baseline (0.73). Ahmed et al. achieved the best F1-score for infection classification (0.68), an improvement of 0.12 on the baseline (0.56). Bloch et al. achieved the best F1-score for ischaemia classification (0.56) which shows an improvement of 0.12 on the baseline (0.44). For F1-score of detection of both infection and ischaemia, Galdran et al. achieved a value of 0.56, resulting in an improvement of 0.09 over the baseline (0.10).

For the micro-average results, Galdran et al. achieved the highest micro-average F1-score (0.68), which is an improvement of 0.05 over the baseline (0.63). They also achieved the highest AUC (0.91), an improvement of 0.04 on the baseline result (0.87).

For the macro-average results, Bloch et al. achieved the highest macro-average precision (0.62) which is an increase of 0.05 over the baseline (0.57). Galdran et al. achieved the best micro-average recall result (0.66), demonstrating an increase of 0.04 over the baseline (0.62). Galdran et al. achieved the highest macro-average F1-score (0.62), which represents an increase of 0.07 over the baseline (0.55). For macro-average AUC, Galdran et al. achieved the best result with (0.89), which is an increase of 0.03 over the baseline result (0.86).

To summarise, the top 3 highest performing entries came from submissions by Galdran et al., Bloch et al. and Ahmed et al. Galdran et al. achieved the highest F1-score for control (0.76), infection and ischaemia (0.57), micro-average F1-score (0.68), micro-average AUC (0.91), and highest macro-average recall

(0.66), macro-average F1-score (0.62) and macro-average AUC (0.89). Bloch et al. achieved the highest F1-score for ischaemia classification (0.56), and macro-average precision (0.62). Ahmed et al. achieved the highest F1-score for infection classification (0.68).

The results from the challenge represent a modest increase in performance metrics when compared to the baseline results (mean = 0.064, standard deviation = 0.035, error = 0.011). Possible reasons for this include significant class imbalance inherent within the dataset and the small size of the sample images. F1-scores for infection cases (0.68) and ischaemic cases (0.56) are significantly lower than control cases (0.76), which is a possible further reflection of the class imbalance within the dataset. Additional curation together with additional inter- and intra-rater reliability measures may help to further enhance the datasets. However, dataset curation is a difficult and time-consuming task, and presents additional challenges in the form of label noise and artefacts which could affect the true accuracy of models trained on our data [34,35].

Table 2. Summary of the top 10 performing networks for DFUC2021, compared to the baseline result. AP = Abdul-prediction, AE = Adrian-ensemble, AR = Adrian-results, AM = Ahmed-moded, LE = Louise-ensemble, LID = Louise-ID, ST = shimmer-test, XT = xie-test.

Method	Metrics									
	Per class F1-score				micro-average		macro-average			
	Control	Infection	Ischaemia	Both	F1	AUC	Precision	Recall	F1	AUC
Baseline [29]	0.73	0.56	0.44	0.47	0.63	0.87	0.57	0.62	0.55	0.86
AP_vit_bas_GP4EVbn	0.74	0.61	0.43	0.33	0.63	0.87	0.52	0.59	0.53	0.85
AP_vit_mil_UNKBe8A	0.73	0.55	0.52	0.42	0.62	0.87	0.57	0.59	0.56	0.84
AP_vit_multi1_test	0.75	0.63	0.47	0.43	0.66	0.87	0.58	0.61	0.57	0.85
AE_bit_effb3_F2	0.76	0.64	0.53	0.56	0.68	0.91	0.61	0.65	0.62	0.89
AE_results_final_test	0.74	0.61	0.49	0.49	0.65	0.90	0.58	0.61	0.58	0.88
AE_results_final_test2	0.76	0.64	0.51	0.54	0.68	0.90	0.61	0.65	0.61	0.88
AR_final_test4	0.73	0.63	0.52	0.50	0.66	0.90	0.59	0.62	0.60	0.87
AR_final_test5	0.75	0.63	0.51	0.57	0.67	0.90	0.61	0.66	0.62	0.88
AM_v0.89_test	0.72	0.67	0.46	0.54	0.67	0.89	0.60	0.60	0.60	0.86
AM_v0.89_test_1	0.71	0.68	0.46	0.53	0.67	0.89	0.60	0.60	0.60	0.86
Arnab	0.73	0.57	0.40	0.45	0.62	0.88	0.53	0.57	0.54	0.85
LE_predictio_aCYsozF	0.75	0.59	0.56	0.54	0.65	0.87	0.62	0.62	0.61	0.86
LE_Predictio_SSufTEW	0.74	0.60	0.52	0.52	0.65	0.89	0.60	0.62	0.59	0.87
LID47_predictions	0.74	0.58	0.54	0.51	0.65	0.87	0.61	0.61	0.60	0.85
LID48_predictions	0.74	0.59	0.55	0.51	0.65	0.85	0.61	0.62	0.60	0.84
LID49_predictions	0.74	0.59	0.54	0.53	0.65	0.87	0.61	0.63	0.60	0.86
Orhun	0.74	0.55	0.52	0.44	0.62	0.89	0.59	0.58	0.56	0.87
ST_submit1	0.73	0.55	0.48	0.41	0.61	0.87	0.57	0.60	0.54	0.86
ST_submit2	0.74	0.60	0.45	0.43	0.64	0.88	0.56	0.59	0.55	0.86
ST_submit3	0.74	0.57	0.49	0.38	0.63	0.87	0.57	0.59	0.55	0.86
ST_submit4	0.73	0.57	0.46	0.46	0.63	0.87	0.57	0.58	0.56	0.86
Weilunwang	0.70	0.54	0.42	0.47	0.60	0.86	0.54	0.57	0.53	0.82
Yeah	0.72	0.55	0.47	0.35	0.60	0.74	0.53	0.56	0.52	0.70
xie-s.9163	0.72	0.52	0.50	0.46	0.60	0.87	0.57	0.58	0.55	0.86
XT_eff_dense_91_I0yNTTP	0.74	0.56	0.46	0.46	0.63	0.88	0.57	0.60	0.55	0.87

3.2 Ensemble of the Top 10 Performing Models

In this section, we analyse the results of ensembling the top performing models submitted to DFUC2021 to determine if an ensemble approach can provide an increase to performance metrics in multi-classification of DFU. Table 3 shows the results of ensembling the top 10 performing models from DFUC2021.

Table 3. Summary of the results for the top 10 teams and further analysis on the ensembled results. Ensemble Top 2 represents an ensemble of the top 2 teams results, Ensemble Top 3 represents an ensemble of the top 3 teams results, etc.

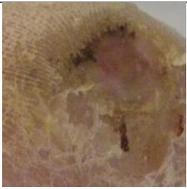
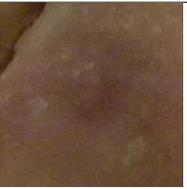
Method	Metrics									
	Per class F1-score				micro-average		macro-average			
	Control	Infection	Ischaemia	Both	F1	AUC	Precision	Recall	F1	AUC
Top-1	0.7574	0.6388	0.5282	0.5619	0.6801	0.9071	0.6140	0.6522	0.6216	0.8855
Top-2	0.7453	0.5917	0.5580	0.5359	0.6532	0.8734	0.6207	0.6246	0.6077	0.8616
Top-3	0.7157	0.6714	0.4574	0.5390	0.6714	0.8935	0.5984	0.5979	0.5959	0.8644
Top-4	0.7466	0.6281	0.4670	0.4347	0.6577	0.8731	0.5814	0.6104	0.5691	0.8488
Top-5	0.7360	0.5468	0.5216	0.4396	0.6199	0.8865	0.5917	0.5759	0.5610	0.8702
Top-6	0.7320	0.5732	0.4621	0.4599	0.6292	0.8725	0.5692	0.5823	0.5568	0.8635
Top-7	0.7407	0.5566	0.4602	0.4558	0.6253	0.8821	0.5705	0.6032	0.5533	0.8698
Top-8	0.7275	0.5701	0.4000	0.4463	0.6222	0.8821	0.5329	0.5692	0.5360	0.8471
Top-9	0.6996	0.5412	0.4237	0.4657	0.5999	0.8622	0.5371	0.5681	0.5326	0.8222
Top-10	0.7192	0.5456	0.4656	0.3532	0.6027	0.7443	0.5300	0.5611	0.5209	0.7020
Ensemble Top 2	0.7455	0.6014	0.5615	0.5301	0.6572	0.9054	0.6187	0.6297	0.6096	0.8866
Ensemble Top 3	0.7491	0.6303	0.5637	0.5799	0.6756	0.9096	0.6352	0.6422	0.6307	0.8870
Ensemble Top 4	0.7578	0.6410	0.5513	0.5412	0.6805	0.9102	0.6244	0.6416	0.6228	0.8893
Ensemble Top 5	0.7571	0.6337	0.5653	0.5411	0.6775	0.9112	0.6314	0.6395	0.6243	0.8933
Ensemble Top 6	0.7566	0.6287	0.5437	0.5383	0.6740	0.9104	0.6253	0.6323	0.6168	0.8947
Ensemble Top 7	0.7611	0.6244	0.5417	0.5137	0.6720	0.9099	0.6201	0.6290	0.6102	0.8968
Ensemble Top 8	0.7618	0.6240	0.5486	0.5309	0.6738	0.9106	0.6251	0.6357	0.6163	0.8980
Ensemble Top 9	0.7615	0.6215	0.5629	0.5292	0.6730	0.9093	0.6291	0.6396	0.6188	0.8964
Ensemble Top 10	0.7628	0.6242	0.5603	0.5209	0.6740	0.9082	0.6280	0.6381	0.6171	0.8954

3.3 Visual Comparison of the Top-10 Methods

We conducted further analysis on the top performing methods to determine trends in the data for images that were both easily predicted correctly with high confidence and images where correct classification was difficult. We then visualised those images and identify key features that could have effected the classification result.

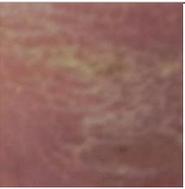
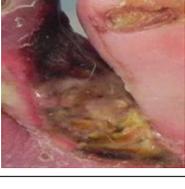
Table 4 highlights 3 images from each class that were correctly classified with high confidence by the top 10 challenge participants. The images show that when the wound is fully visible, the networks are able to determine the difference between classes, even when features from other classes are present, e.g. None image 1 where the black section could cause Ischaemia bias. In-contrast 5 show examples of images that were incorrectly classified by all top 10 challenge participants. These examples highlight the issue with extreme angles and image blur, as seen in images 1 and 2 for the None class, and in image 2 for Both.

Table 4. Images from the testing set which the top 10 networks all predicted correctly.

Class	Image 1	Image 2	Image 3
None			
Infection			
Ischaemia			
Both			

One particularly notable result can be seen when comparing the correctly classified image in Table 4 (Image 1, Both) with the incorrectly classified image in Table 5 (Image 1, Both). The image in Table 5 is the result of subtle natural augmentation and has resulted in an incorrect classification.

Table 5. Images from the testing set which the top 10 networks all predicted incorrectly.

Class	Image 1	Image 2	Image 3
None			
Infection			
Ischaemia			
Both			

4 Conclusion

In this study, we introduce the largest DFU pathology dataset, and propose a weakly supervised framework for DFU pathology classification of infection and ischaemia. This is the first dataset of its kind to be made available to the research community together with implementation of CNNs for multi-class classification of infection, ischaemia, and co-occurrences of infection and ischaemia. These advancements will help to support early identification of DFU complications to

guide treatment and help prevent further complications including limb amputation.

Although the majority of the deep learning methods reported in this paper show promising results in recognising infection and ischaemia, there are still significant challenges in designing methods to detect the co-occurrences of both conditions. Future work will investigate more advanced techniques such as generative adversarial networks and unsupervised learning to improve network performance.

This work will form an important contribution to our ongoing research into developing a fully automated DFU diagnosis and monitoring framework which can be used by patients and their carers in home settings, to help reduce strains on healthcare services around the world. This work will build on our existing framework [36,37] in delivering an easy-to-use system capable of advanced forms of diabetic foot analysis, which will include longitudinal monitoring as a means of assessing wound healing progress.

Acknowledgment

We gratefully acknowledge the support of NVIDIA Corporation who provided access to GPU resources for the DFUC2021 Challenge and an NVIDIA Geforce RTX 3090 GPU card as the prize for the winning team.

References

1. David G Armstrong, Andrew JM Boulton, and Sicco A Bus. Diabetic foot ulcers and their recurrence. *New England Journal of Medicine*, 376(24):2367–2375, 2017.
2. Andrew James Michael Boulton, David G Armstrong, Robert S Kirsner, Christopher E Attinger, Lawrence A Lavery, Benjamin A Lipsky, Joseph L Mills Sr, and John S Steinberg. Diagnosis and management of diabetic foot complications. 2019.
3. Changhan Wang, Xinchun Yan, Max Smith, Kanika Kochhar, Marcie Rubin, Stephen M Warren, et al. A unified framework for automatic wound segmentation and analysis with deep convolutional neural networks. In *Engineering in Medicine and Biology Society (EMBC), 2015 37th Annual International Conference of the IEEE*, pages 2415–2418. IEEE, 2015.
4. M. Goyal, M. H. Yap, N. D. Reeves, S. Rajbhandari, and J. Spragg. Fully convolutional networks for diabetic foot ulcer segmentation. In *2017 IEEE International Conference on Systems, Man, and Cybernetics (SMC)*, pages 618–623, Oct 2017.
5. M. Goyal, N. D. Reeves, A. K. Davison, S. Rajbhandari, J. Spragg, and M. H. Yap. Dfunet: convolutional neural networks for diabetic foot ulcer classification. *IEEE Transactions on Emerging Topics in Computational Intelligence*, pages 1–12, 2018.
6. M. Goyal, N. D. Reeves, S. Rajbhandari, and M. H. Yap. Robust methods for real-time diabetic foot ulcer detection and localization on mobile devices. *IEEE Journal of Biomedical and Health Informatics*, 23(4):1730–1741, July 2019.
7. Moi Hoon Yap, Ryo Hachiuma, Azadeh Alavi, Raphael Brüngel, Bill Cassidy, Manu Goyal, Hongtao Zhu, Johannes Rückert, Moshe Olshansky, Xiao Huang, Hideo Saito, Saeed Hassanpour, Christoph M. Friedrich, David B. Ascher, Anping Song, Hiroki Kajita, David Gillespie, Neil D. Reeves, Joseph M. Pappachan,

- Claire O'Shea, and Eibe Frank. Deep learning in diabetic foot ulcers detection: A comprehensive evaluation. *Computers in Biology and Medicine*, 135:104596, 2021.
8. F William Wagner. The diabetic foot. *Orthopedics*, 10(1):163–172, 1987.
 9. Lawrence A Lavery, David G Armstrong, and Lawrence B Harkless. Classification of diabetic foot wounds. *The Journal of Foot and Ankle Surgery*, 35(6):528–531, 1996.
 10. David G Armstrong, Lawrence A Lavery, and Lawrence B Harkless. Validation of a diabetic wound classification system: the contribution of depth, infection, and ischemia to risk of amputation. *Diabetes care*, 21(5):855–859, 1998.
 11. Paul Ince, Zulfiqarali G Abbas, Janet K Lutale, Abdul Basit, Syed Mansoor Ali, Farooq Chohan, Stephan Morbach, et al. Use of the sinbad classification system and score in comparing outcome of foot ulcer management on three continents. *Diabetes care*, 31(5):964–967, 2008.
 12. Joseph L Mills Sr, Michael S Conte, David G Armstrong, Frank B Pomposelli, Andres Schanzer, Anton N Sidawy, et al. The society for vascular surgery lower extremity threatened limb classification system: risk stratification based on wound, ischemia, and foot infection (wif). *Journal of vascular surgery*, 59(1):220–234, 2014.
 13. David G Armstrong and Joseph L Mills. Juggling risk to reduce amputations: the three-ring circus of infection, ischemia and tissue loss-dominant conditions. *Wound Medicine*, 1:13–14, 2013.
 14. Maximiano Albers, Ayrton C Fratezi, and Nelson De Luccia. Assessment of quality of life of patients with severe ischemia as a result of infrainguinal arterial occlusive disease. *Journal of vascular surgery*, 16(1):54–59, 1992.
 15. L Prompers, M Huijberts, Jan Apelqvist, E Jude, A Piaggese, K Bakker, et al. High prevalence of ischaemia, infection and serious comorbidity in patients with diabetic foot disease in europe. baseline results from the eurodiale study. *Diabetologia*, 50(1):18–25, 2007.
 16. Benjamin A Lipsky, Anthony R Berendt, Paul B Cornia, James C Pile, Edgar JG Peters, David G Armstrong, et al. 2012 infectious diseases society of america clinical practice guideline for the diagnosis and treatment of diabetic foot infections. *Clinical Infectious Diseases*, 54(12):e132–e173, 2012.
 17. Lawrence A Lavery, David G Armstrong, Robert P Wunderlich, Jeffrey Tredwell, and Andrew JM Boulton. Diabetic foot syndrome: evaluating the prevalence and incidence of foot pathology in mexican americans and non-hispanic whites from a diabetes disease management cohort. *Diabetes Care*, 26(5):1435–1438, 2003.
 18. Grant H Skrepnek, Joseph L Mills, Lawrence A Lavery, and David G Armstrong. Health care service and outcomes among an estimated 6.7 million ambulatory care diabetic foot cases in the us. *Diabetes Care*, 40(7):936–942, 2017.
 19. Jaap J van Netten, Damien Clark, Peter A Lazzarini, Monika Janda, and Lloyd F Reed. The validity and reliability of remote diabetic foot ulcer assessment using mobile phone images. *Scientific Reports*, 7(1):9480, 2017.
 20. Mark Swerdlow, Laura Shin, Karen D'Huyvetter, Wendy J. Mack, and David G. Armstrong. Initial clinical experience with a simple, home system for early detection and monitoring of diabetic foot ulcers: The foot selfie. *Journal of Diabetes Science and Technology*, 2021.
 21. Moi Hoon Yap, Katie E Chatwin, Choon-Ching Ng, Caroline A Abbott, Frank L Bowling, Satyan Rajbhandari, et al. A new mobile application for standardizing diabetic foot images. *Journal of diabetes science and technology*, 12(1):169–173, 2018.

22. Bill Cassidy, Neil D Reeves, Joseph M Pappachan, David Gillespie, Claire O'Shea, Satyan Rajbhandari, Arun G Maiya, Eibe Frank, Andrew J M Boulton, David G Armstrong, Bijan Najafi, Justina Wu, Rupinder Singh Kochhar, and Moi Hoon Yap. The dfuc 2020 dataset: Analysis towards diabetic foot ulcer detection. *touchREVIEWS in Endocrinology*, 17:5–11, 2021.
23. Andre Esteva, Brett Kuprel, Roberto Novoa, Justin Ko, Susan Swetter, Helen Blau, and Sebastian Thrun. Dermatologist-level classification of skin cancer with deep neural networks. *Nature*, 542, 01 2017.
24. Titus J. Brinker, Achim Hekler, Alexander H. Enk, Carola Berking, Sebastian Haferkamp, Axel Hauschild, Michael Weichenthal, Joachim Klode, Dirk Schadendorf, Tim Holland-Letz, Christof von Kalle, Stefan Fröhling, Bastian Schilling, and Jochen S. Utikal. Deep neural networks are superior to dermatologists in melanoma image classification. *European Journal of Cancer*, 119:11–17, 2019.
25. Y. Fujisawa, Y. Otomo, Y. Ogata, Y. Nakamura, R. Fujita, Y. Ishitsuka, R. Watanabe, N. Okiyama, K. Ohara, and M. Fujimoto. Deep-learning-based, computer-aided classifier developed with a small dataset of clinical images surpasses board-certified dermatologists in skin tumour diagnosis. *The British Journal of Dermatology*, 180(2):373–381, 2019.
26. T. C. Pham, V. D. Hoang, C. T. Tran, M. S. K. Luu, D. A. Mai, A. Doucet, and C. M. Luong. Improving binary skin cancer classification based on best model selection method combined with optimizing full connected layers of deep cnn. In *2020 International Conference on Multimedia Analysis and Pattern Recognition (MAPR)*, pages 1–6, 2020.
27. Shunichi Jinnai, Naoya Yamazaki, Yuichiro Hirano, Yohei Sugawara, Yuichiro Ohe, and Ryuji Hamamoto. The development of a skin cancer classification system for pigmented skin lesions using deep learning. *Biomolecules*, 10(8), 2020.
28. Manu Goyal, Neil D. Reeves, Satyan Rajbhandari, Naseer Ahmad, Chuan Wang, and Moi Hoon Yap. Recognition of ischaemia and infection in diabetic foot ulcers: Dataset and techniques. *Computers in Biology and Medicine*, 117:103616, 2020.
29. Moi Hoon Yap, Bill Cassidy, Joseph M Pappachan, Claire O'Shea, David Gillespie, and Neil Reeves. Analysis towards classification of infection and ischaemia of diabetic foot ulcers. *arXiv preprint arXiv:2104.03068*, 2021.
30. George Forman and Martin Scholz. Apples-to-apples in cross-validation studies: pitfalls in classifier performance measurement. *Acm Sigkdd Explorations Newsletter*, 12(1):49–57, 2010.
31. Jingyi Qu, Ting Zhao, Meng Ye, Jiayi Li, and Chao Liu. Flight delay prediction using deep convolutional neural network based on fusion of meteorological data. *Neural Processing Letters*, 52, 10 2020.
32. Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *CoRR*, abs/1409.1556, 2015.
33. Pierre Foret, Ariel Kleiner, Hossein Mobahi, and Behnam Neyshabur. Sharpness-aware minimization for efficiently improving generalization. *arXiv preprint arXiv:2010.01412*, 2021.
34. David Wen, Saad M Khan, Antonio Ji Xu, Hussein Ibrahim, Luke Smith, Jose Caballero, Luis Zepeda, Carlos de Blas Perez, Alastair K Denniston, Xiaoxuan Liu, and Rubeta N Matin. Characteristics of publicly available skin cancer image datasets: a systematic review. *The Lancet Digital Health*, 2021.
35. Bill Cassidy, Connah Kendrick, Andrzej Brodzicki, Joanna Jaworek-Korjakowska, and Moi Hoon Yap. Analysis of the isic image datasets: Usage, benchmarks and recommendations. *Medical Image Analysis*, 2021.

36. Neil D. Reeves, Bill Cassidy, Caroline A. Abbott, and Moi Hoon Yap. Chapter 7 - novel technologies for detection and prevention of diabetic foot ulcers. In Amit Gefen, editor, *The Science, Etiology and Mechanobiology of Diabetes and its Complications*, pages 107–122. Academic Press, 2021.
37. Bill Cassidy, Neil D. Reeves, Joseph M. Pappachan, Naseer Ahmad, Samantha Haycocks, David Gillespie, and Moi Hoon Yap. A cloud-based deep learning framework for remote detection of diabetic foot ulcers. *arXiv preprint arXiv:2004.11853*, 2021.