

# Translating Clinical Delineation of Diabetic Foot Ulcers into Machine Interpretable Segmentation

Connah Kendrick, Bill Cassidy, Joseph M. Pappachan, Claire O'Shea, Cornelious J. Fernandez, Elias Chacko, Koshy Jacob, Neil D. Reeves, and Moi Hoon Yap, *Senior Member, IEEE*

**Abstract**—Diabetic foot ulcer is a severe condition that requires close monitoring and management. For training machine learning methods to auto-delineate the ulcer, clinical staff must provide ground truth annotations. In this paper, we propose a new diabetic foot ulcer dataset, namely DFUC2022, the largest segmentation dataset where ulcer regions were manually delineated by clinicians. We assess whether the clinical delineations are machine interpretable by deep learning networks or if image processing refined contour should be used. By providing benchmark results using a selection of popular deep learning algorithms, we draw new insights into the limitations of DFU wound delineation and report on the associated issues. With in depth understanding and observation on baseline models, we propose a new strategy for training and modify the FCN32 VGG network to address the issues. We achieved notable improvement with a Dice score of 0.7446, when compared to the best baseline network of 0.5708 and the first place in DFUC2022 challenge leaderboard, with a Dice score of 0.7287. This paper demonstrates that image processing using refined contour as ground truth can provide better agreement with machine predicted results. Furthermore, we propose a new strategy to address the limitations of the existing training protocol. For reproducibility, all source code will be made available upon acceptance of this paper, and the dataset is available upon request.

**Index Terms**—Clinical delineation, deep learning, DFUC2022, diabetic foot ulcers, segmentation

## I. INTRODUCTION

**D**IABETIC Foot Ulcers (DFU) are caused when sections of the foot and skin are damaged due to multiple factors including nerve damage (diabetic peripheral neuropathy) and foot deformities. DFU healing can be impaired due to blood flow (vascular) limitations as a consequence of diabetes. Owing to this, the DFU requires regular checks to ensure optimal healing and to inform any adjustments to the treatment strategy. DFU frequently become infected, can lead to

amputation and in some cases loss of life if antibiotic treatment is unsuccessful [1].

It is shown that at least 10% of people with diabetes will have some form of DFU in their lifetime, rising to 25% depending on life-style factors [2], [3]. Moreover, recent studies have shown that after treatment, patients have a 70% chance of ulcer recurrence [4]. Although DFU is a physical disease, it has also been widely reported to have a drastic impact on patient mental well-being and quality of life, causing anxiety and depression [5].

Treatment for DFU can be a long-term process, due to diabetes-related complications impairing the healing process [6]. It requires a multi-disciplinary team [7] to monitor the progress of the ulcer, focusing largely on the management of diabetes [8] and blood flow to the foot. However, complications, such as infection [9] significantly prolong treatment. If treatment is prolonged, the possibility of infection and amputation increase significantly [10]. This has been shown to create a heavy burden on healthcare systems, in terms of both time and cost per patient [7], [11]. Furthermore, this causes a great deal of concern due to the predicted rapid global rise of diabetes [12], amplified significantly by the current pandemic [13]. To address these challenges, researchers have been working towards development of methods [14]–[18] and automated systems capable of detecting and monitoring DFU [19], [20]. Improvements to automated delineation of DFU could support improved digital healthcare tools that could be used for screening and triage of DFU. Furthermore, these improvements could aid in the development of active DFU monitoring systems, to engage the healing process stage.

This paper demonstrates the processes of translating clinical delineation of DFU into machine interpretable segmentation. We contribute to the research progress of DFU segmentation in the following ways:

- Introduce the largest DFU segmentation dataset to date with ground truth delineation (namely, DFUC2022) and perform detailed analysis.
- Investigate the effect of image processing refined contours on the performance of a popular deep learning segmentation algorithm, DeepLabv3+.
- Establish baseline results for the DFUC2022 dataset using a range of popular deep learning segmentation networks.
- Propose a new strategy to optimise the performance of DFU segmentation in an end-to-end network and achieved the best result when compared to the DFUC2022 challenge leaderboard's results.

Date submitted for review: 30 September 2022. We gratefully acknowledge the support of NVIDIA Corporation who provided access to GPU resources and sponsored Diabetic Foot Ulcers Grand Challenges.

C. Kendrick, B. Cassidy and M.H. Yap are with the Department of Computing and Mathematics, Manchester Metropolitan University (e-mail: Connah.Kendrick@mmu.ac.uk, M.Yap@mmu.ac.uk).

J. Pappachan is with the Lancashire Teaching Hospitals NHS Foundation Trust (e-mail: pappachan.joseph@lthtr.nhs.uk).

C. O'Shea is with Waikato District Health Board (e-mail: claire.o'shea@waikatodhb.health.nz).

C. J. Fernandez is with United Lincolnshire Hospitals NHS Trust (e-mail: drcjfernandez@yahoo.com).

E. Chacko is with Jersey General Hospital (e-mail: e.chacko@health.gov.je).

K. Jacob is with Eastbourne District General Hospital (e-mail: drkoshyjacob@gmail.com).

N. D. Reeves is with the Research Centre for Musculoskeletal Science & Sports Medicine (e-mail: n.reeves@mmu.ac.uk).

TABLE I  
A COMPARISON OF THE PROPOSED DFUC2022 DATASETS AND THE EXISTING DFU IMAGE SEGMENTATION DATASETS.

Publication	Year	Dataset Name	Resolution	Train	Test	Total
Wang et al. [21]	2020	AZH wound care dataset	$224 \times 224$	831	278	1109
Thomas [22]	NA	Medetec	$560 \times 391$ $224 \times 224$	152	8	160
Wang et al. [23]	2021	FUSeg Challenge	$512 \times 512$	1010	200	1210
Proposed	2022	DFUC2022	$640 \times 480$	2000	2000	4000

This work will benefit the research community by providing a summary of available datasets to access and use for training segmentation based networks. With our established partnerships between clinicians and researchers, we provide the largest DFU segmentation dataset with superior image resolution when compared with existing DFU datasets [23]. Additionally, we provide an in-depth analysis on the performance of baseline results and propose a new end-to-end network, resulting in superior performance when compared to the best reported model in the challenge leaderboard. To assist in fair assessment and comparison with the benchmarks, we release a testing set that can be evaluated online via a grand challenge website, providing almost instant evaluation results on a standard set of performance metrics.

## II. RELATED WORK

### A. Previous Datasets

1) *DFUC2020 Dataset*: The DFUC2020 Dataset [14] is an object detection based dataset, containing 2000 training, 200 validation and 2000 testing images. All images are  $640 \times 480$ , but some images contained multiple DFUs, increasing the total number of detection annotations. Three cameras were used for image capture, Kodak DX4530, Nikon D3300 and Nikon COOLPIX P100. The images were acquired with close-ups of the full foot at a distance of around 30–40 cm with the parallel orientation to the ulcer. The use of flash was avoided, and instead, room lights were used to provide consistent colours in the images. Images were acquired by a podiatrist and a consultant physician with specialization in the diabetic foot, both with more than 5 years professional experience. All images were captured without the use of a tripod.

2) *DFUC2021 Dataset*: The DFUC2021 dataset [15] is a multi-class DFU dataset, targeting DFU, infection, Ischaemia and both. The dataset contains 5,955 training images, and 5,734 for testing. Additionally, 3,994 images were released unlabeled to support semi and self-supervised methods. Images were captured under the same setting as the DFUC2020 dataset.

3) *FUSeg dataset*: Wang et al. [24] introduced the Foot Ulcer Segmentation Dataset. This work focused on the development of segmentation CNNs using 1210 foot photographs exhibiting DFU which were collected over a 2 year period from 889 patients. They provided ground truth masks provided by wound care experts. However, many of the images were heavily padded to standardise image dimensions for training purposes. Additionally, although the images were shared as lossless PNG files, they exhibited notable compression artefacts, indicating that the original images had been heavily

compressed before being converted to PNG. The provided ground truth files also appeared to be a mix of human and machine-generated masks. The images were  $512 \times 512$  with 1000 for training and 200 for test. The capture equipment was a Canon SX 620 hs and an iPad Pro. The AZH wound care and Medetec datasets, see Table I, where both used as part of the FUSeg dataset. It is noted that the AZH dataset is cropped to the ulcer region, where as the final images in the FUSeg challenge have surrounding regions.

### B. Related Methods

The first works in DFU segmentation using fully convolutional techniques were completed by Goyal et al. [25]. They performed segmentation experiments using a small dataset comprising 705 images with an FCN-16s network. They used 5-fold cross-validation with two-tier transfer learning resulting in a Dice Similarity Coefficient of 0.794 ( $\pm 0.104$ ) for segmentation of DFU region. These results were promising, however, the small size of the dataset is likely to impact the model’s ability to generalise in real-world use.

More recently, the winning team of the FUSeg challenge, Mahbod et al. [26], used an ensemble of LinkNet and U-Net networks. They achieved a Dice Similarity Coefficient of 0.888. They used pretrained weights (EfficientNetB1 for LinkNet and EfficientNetB2 for U-Net) with additional pre-training using the Medetec dataset. The challenge concluded that segmentation of small isolated areas of the wound with ambiguous boundaries were the most challenging aspects of the task. Conversely, segmentation of relatively larger wound regions showing clear boundaries where wound beds were cleansed, removing dead tissue, provided superior results. Cases clearly exhibiting infection, slough, or other impediments were also shown to provide improved results.

Current works in DFU segmentation show promising results. However, there are notable limitations to the datasets that were used to train these models. Aspects such as the quality and number of images may present issues that would negatively affect real-world application.

## III. THE DFUC2022 DATASET

We have received approval from the UK National Health Service Research Ethics Committee (reference number is 15/NW/0539) to use DFU images for the purpose of this research. This paper introduces the largest DFU segmentation dataset, which consists of a training set of 2000 images and a testing set of 2000 images.

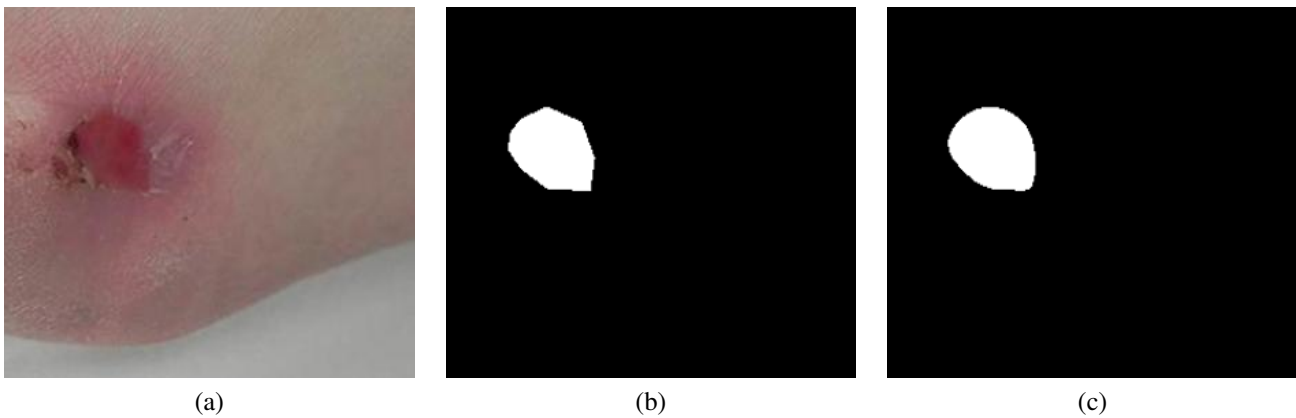


Fig. 1. Illustration of (a) an early onset DFU; (b) expert delineation; and (c) refined active contour shape.

### A. Dataset Construction

The dataset was constructed in collaboration with the medical experts from Lancashire Teaching Hospitals, Waikato District Health Board, United Lincolnshire Hospitals, Jersey General Hospital, and Eastbourne District General Hospital. The DFUs were captured at room lighting, in full foot view, around 30-40cm away with the DFU centered. Three cameras were used, i.e., Kodak DX4530, Nikon D3300 and Nikon COOLPIX P100. All images were taken by experienced podiatrist and physician in foot clinic. Images were then downsampled to  $640 \times 480$  and stored in JPG format.

### B. Reference Annotation Protocol

The ulcer regions on these images were delineated by experienced podiatrists and consultants. The podiatrists used the VGG annotator software, to produce a polygon outline of the DFU region in JSON format. The JSON files were then converted into binary mask images and stored in PNG format. We then preprocess the raw masks with an active contour algorithm [27].

Figure 1 illustrates an example of a DFU image showing a preprocessed region with active contour together with the expert delineation. Note that the boundary of the region is smoother after the preprocessing stage. To ensure that this smoothing process does not alter the clinical delineation, we report the agreement between expert delineation and refined contours, which produced a high agreement rate with a Dice Score of  $0.9650 \pm 0.0226$  and Mean Intersection Over Union (mIoU) of  $0.9332 \pm 0.0408$ . These metrics demonstrate that preprocessing did not significantly alter clinical delineation, where the number of DFUs are equivalent before and after preprocessing.

The DFUC2022 training set consists of 2304 ulcers, where the smallest ulcer size is 0.04% of the total image size, and the largest ulcer size is 35.04% of the total image size. Figure 2 provides an overview of the ratio of the delineated ulcer region to the total image size, where 89% (2054 out of 2304) of the ulcers are less than 5% of the total image size. The smaller images in particular represent a significant challenge for segmentation algorithms as it is widely known that deep learning algorithms have a tendency to miss small regions [28].

Another advantage of our dataset is that of the 2000 training images, there are 2304 ulcers with an average of 1.152 ulcers per image.

## IV. METHODS

This section describes the methods used to investigate the effect of image processing refined contours, summarises a range of popular baseline methods for medical image segmentation, and a new strategy to improve the performance of the best segmentation method on the DFUC2022 dataset. We provide segmentation masks for the training set only, and use the grand-challenge website (<https://dfuc2022.grand-challenge.org/>) to allow researchers to test their methods on an exclusive testing set. We provide a total of 4000 images with 2000 binary masks for training. The masks are coded 0 for background and 1 for the DFU region.

### A. Manual delineation vs refined contours

While deep learning has gained popularity in biomedical image segmentation, there are unanswered questions concerning ground truth annotation, such as: (1) would deep learning algorithms learn better with expert manual delineations (polygonal outlines) or image processing refined contours; and (2) which contour should be used for machine learning algorithms? To answer these questions in the context of DFUC2022, we run experiments with Deeplabv3+ [29], one of the popular deep learning algorithms for medical imaging research [30], [31]. Our intention is not to produce the best result, but to study the effect of coarse and detailed delineation on deep learning algorithms. Therefore, we select this algorithm without bias. First, we train two models using the default setting of Deeplabv3+, one on expert delineation and another on refined contour. We split the 2000 training images into 1800 images as training set and 200 images as validation set. Then, we test each model on the 2000 test set by using both expert delineation and refined contour as ground truth.

### B. Baseline methods

We implement a wide range of existing deep learning segmentation models for the DFUC2022 baseline. These models cover a range of segmentation architectures, namely FCN

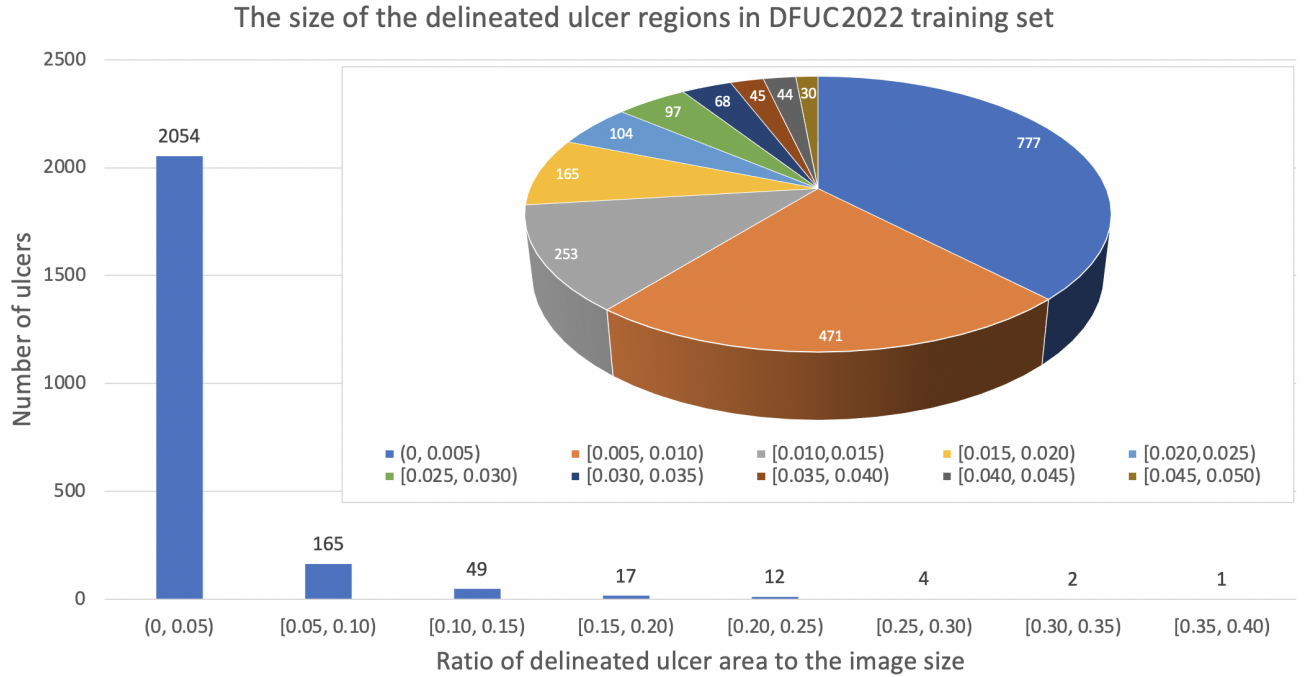


Fig. 2. The size distribution of delineated ulcer regions in the DFUC2022 training set. It is noted that the majority of the ulcers are smaller than 5% of the image size.

[32], U-Net [33] and SegNet [34] with varying backbones to process the data, such as VGG [35] and ResNet50 [36]. We also include a comparison of alternative network depths. The range of model diversity aims to provide a good indication of techniques suitable for DFU segmentation. These new insights can direct future works with a baseline to compare against and reduce the need for repeat training of these networks. In addition to the standard U-Net, SegNet models, we provide baselines for FCN8, FCN32, U-Net and SegNet with ResNet50 and VGG as backbones.

For training the baseline networks, we use all 2000 training images, with 200 separate images for validation. We train the networks with the AdaDelta optimizer and a suggested learning rate of 0.001, decay of 0.95, a stabilisation epsilon of  $1e - 07$  as illustrated in Equation 1, and using categorical cross-entropy loss, as in Equation 2.

$$E[g^2]_t = \rho E[g^2]_{t-1} + (1 - \rho)g_t^2 \quad (1)$$

$$CE = - \sum_{i=1}^o Y_i \cdot \log \hat{X}_i \quad (2)$$

where  $Y_i$  is the  $i$ -th ground truth value and  $\hat{X}_i$  is the predicted value  $a_i$ . We train on multiple batch sizes (2, 32, and 96), Equation 3 and report the best result, as defined by [37].

$$\Delta w_t = - \frac{\eta}{\sqrt{E[g^2]_t + \varepsilon}} \quad (3)$$

We do not perform augmentation during training or post-processing on the final prediction masks, as our aim is to produce baselines and understanding of the DFUC2022 dataset. We train the networks until the validation accuracy fails to improve, with a patience of 10 epochs.

### C. Challenge Competition

To enable open research the DFUC2022 dataset was released in three parts between the 27th April 2022 and the 1st July 2022:

- Training dataset, 2000 images: 27th April 2022.
- Validation dataset, 200 images: 21st June 2022.
- Test dataset, 2000 images: 1st July 2022.

At the release of the validation and test dataset, we released online submissions for live testing. We closed the online submissions on the 29th July 2022, during this time participants could analyse their methods via the validation scores. After the release of test results, we opened a live testing leaderboard to allow future submissions. We compare against the top-10 results in the challenge leaderboard.

### D. Proposed method

Results from the baseline models highlighted a number of issues, such as pixelation and a high number of false positives (small regions). Previous research uses post-processing methods to improve performance. Instead of using morphology, we propose a new strategy using an modified end-to-end deep learning network to enable improved learning of our dataset, and remove the post-processing process. We use the FCN32 architecture with VGG as backbone, as shown in Figure 3. First, we replace the standard ReLU layer in the full network with Leaky-ReLU, depicted by Equation 4.

$$f(x) = \begin{cases} \alpha x & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases} \quad (4)$$

where  $\alpha$  is a scalar for sub zero values and  $x$  the input, with an alpha of 0.3 which aids network learning as it prevents

TABLE II  
COMPARISON OF THE RESULTS WHEN TRAINING WITH MANUAL DELINEATION AS GROUND TRUTH VS IMAGE PROCESSING REFINED CONTOUR AS GROUND TRUTH. THE RESULTS SHOW THE MACHINE PREDICTED MASKS HAVE BETTER AGREEMENT WITH THE REFINED CONTOUR.

Train	Test	Metrics	
		Dice	mIoU
Manual delineation	Manual delineation	0.5870±0.3135	0.4809±0.2993
Manual delineation	Refined contour	0.5930±0.3131	0.4871±0.2999
Refined contour	Manual delineation	0.6219±0.0286	0.5162±0.2967
Refined contour	Refined contour	<b>0.6277±0.3051</b>	<b>0.5224±0.2967</b>

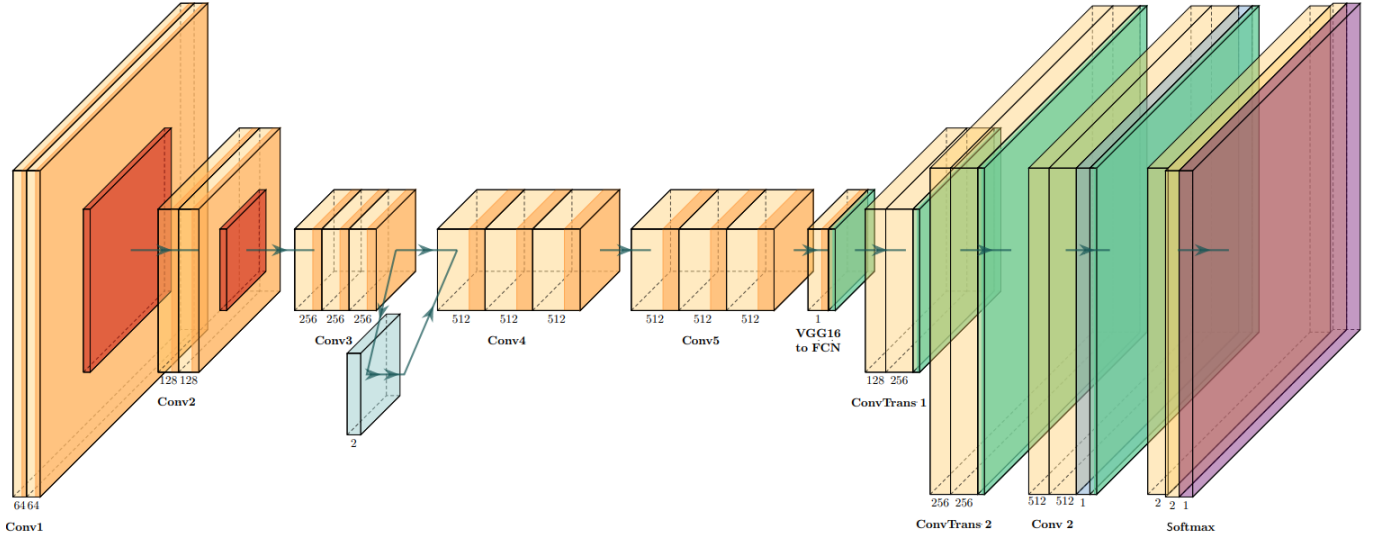


Fig. 3. Illustration of the network structure. Orange represents convolutional layers with Leaky ReLU activation, red indicates max pooling, and light green indicates skip connections using modified squeeze and excite. In the decoding section, green is a dropout layer, yellow is a separable convolution, with dilation, and the softmax layer.

dead neurons occurring. Then, we target excessive down-sampling by removing the bottom three max-pooling layers, while maintaining the padding. This process retains the feature map size from the lowest in the standard network of  $20 \times 15$  to  $160 \times 120$  on the full size images, improving the ability of the network to maintain feature maps of smaller ulcers and tracks overall wound shape, which reduces the issues with biases in dataset distribution.

To resolve the issue of background noise, we experimented using gated convolutions [38]. During this stage the best performing method was modified using a squeeze and excite layer [39] after the final pooling was used, where a dilated convolution (kernel size 5, dilation rate  $2 \times 2$ ) focused on separating the foot region features from the background, which had a standard convolution (kernel size  $1 \times 1$ ) with sigmoid activation. The resulting feature map was multiplied against the normal output of the 3rd pooling stage of the network. These adjustments resulted in improved removal of noisy inconsistent data, reducing the background features of the environment and improving focus on the more consistent foot regions. Thus, the lower levels of the networks can separate the similar textural features of the DFU and foot region. We then address the issue of rapid up-sampling by adjusting the FCN network to gradually grow the predictions through a series of small transposed convolutions (kernel size  $2 \times 2$ , stride  $2 \times 2$ )

with a convolution to refine the contours of the up-sample until the desired size is reached. In many segmentation tasks, post processing of outputs is performed to smooth predictions and blob removal, however we accomplish this internally within the network with a final dilated convolution, as shown in Equation 5 (kernel size  $3 \times 3$ , dilation rate 2).

$$(F *_l k)(p) = \sum_{s+lt=p} F(s)K(t) \quad (5)$$

where  $l$  is the dilation rate providing a gap between receptive points.  $K(t)$  is the values of the filter.  $F(s)$  is the input to the layer and  $\sum_{s+lt=p}$  is the sum of the receptive fields. This allows for the surrounding regions to determine if the section is a small island for removal, or an edge for smoothing, using the wider receptive field.

We also adjusted the training routine using a weighted loss function, which showed further improvements. However, for better results we used the standard loss function and fully balanced the dataset, we processed the training dataset to crop out sliding windows of  $64 \times 48$  with a stride of  $32 \times 24$ , as illustrated in Figure 4. The stride allowed the network to obtain as much of the wound features as possible, producing a total of 810,000 patches. Next, all the patches from the set that contained no DFU pixels were removed, leaving 55,760

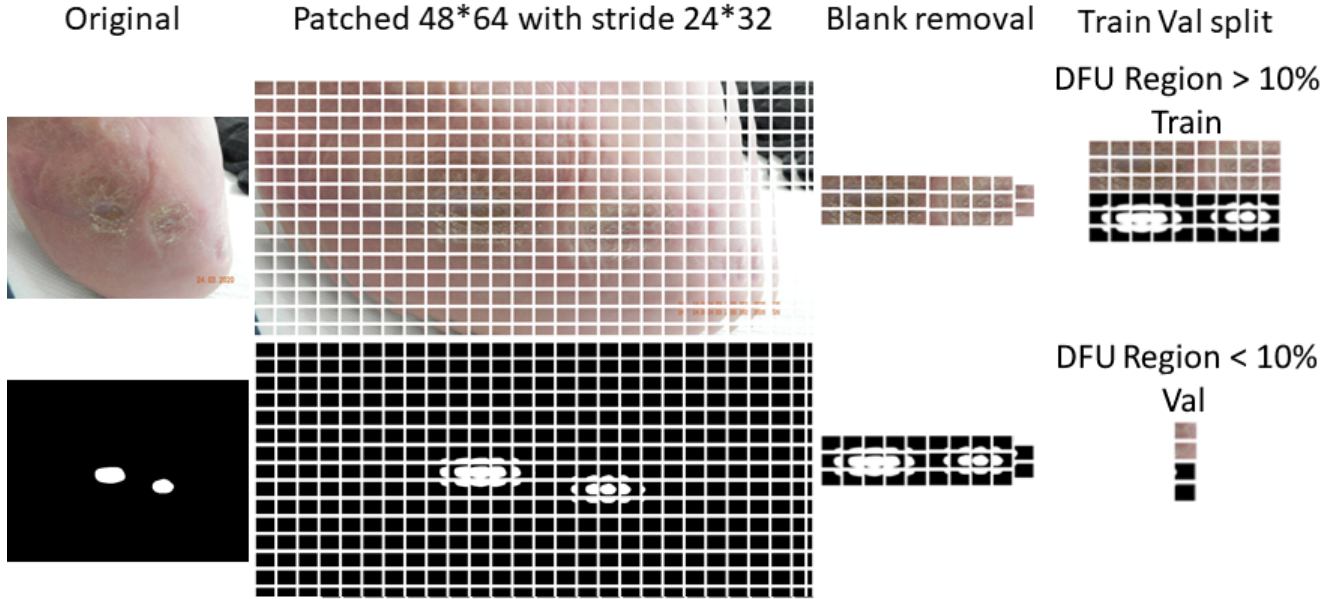


Fig. 4. Illustration of the patching system used for creating the training and validation sets. We use a half stride to create the image windows, to increase the dataset size and reduce chances of only edge cases. We then remove all blank patches from the set, use all with greater than 10% DFU pixels for training and any other for validation.

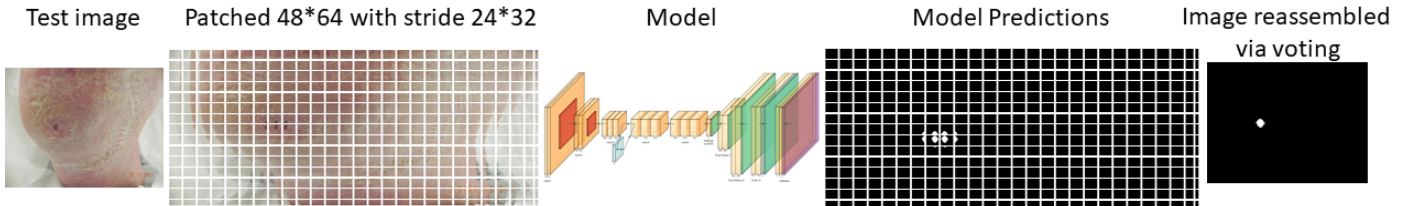


Fig. 5. Illustration of the Testing pipeline. We test on all patches of the images without removal to ensure the system is capable of predicting under a wide array of backgrounds. We then use a vote when reconstructing the image, due to the stride overlap where all must agree to be classed as an ulcer.

patches with DFU pixels. After this, we processed the images to create the training and validation sets, by moving any images with less than 10% DFU pixels into the validation set and using all others for training, giving a total of 38,997 patches for training and 16,763 patches for validation. This stage provided two key advantages:

- *Balanced split of classes:* In total the amount of background pixels was 51.71% and DFU pixels was 48.29%. Thus, giving a more balanced set compared to the standard training method, meaning that both classes will have even weighting.
- *Difficult validation set:* The validation set was heavily biased towards background features. Many of the validation case were small edge cases which are particularly challenging for segmentation networks. This means that a good score reflects a network with clear data understanding.

For the modified network, we train on a batch size of 2, providing the network a balanced view of the data. The same settings for optimizer, learning rate, and loss function are used as in the baseline methods. The network was adjusted to take in the patches at their current resolution. For the modified network, the test dataset was also split in the same process

of  $64 \times 48$  with a stride of  $32 \times 24$ . To reconstruct the image overlapping sections, due to the stride, all patches had to agree for the pixel to be classed as ulcer, as show in Figure 5.

### E. Performance metrics

In image segmentation, the commonly used evaluation metrics are:

Dice Similarity Index as shown in Equation (6):

$$Dice = 2 * \frac{|X \cap Y|}{|X| + |Y| - |X \cap Y|} \quad (6)$$

and Intersection Over Union (IoU) (also known as Jaccard Index) as shown in Equation (7):

$$IoU = \frac{|X \cap Y|}{|X| + |Y|} \quad (7)$$

where  $X$  and  $Y$  represent the ground truth mask and the predicted mask. We used mIoU to better represent the segmentation outcomes for both classes (ulcer and background).

We include additional metrics to understand Type I and Type II errors of the algorithm performance. These two additional

TABLE III

A COMPARISON OF THE OVERALL PERFORMANCE OF THE STATE-OF-THE-ART METHODS WITH AND WITHOUT PRETRAINED MODEL, RESULTS REPORTED ON THEIR BEST EPOCH. † = HIGHER SCORE IS BETTER; ‡ = LOWER SCORE IS BETTER. WE TRAIN 12 UNIQUE MODELS WITH DIFFERENT BATCH SIZES. HOWEVER, WE ONLY SHOW THE MODELS WITH THE SETTINGS THAT RESULTED IN THE BEST PERFORMANCE. *Italic* INDICATES THE BEST BASELINE RESULT AND **BOLD** INDICATES THE BEST OVERALL RESULT.

Model	Backbone	Settings	Metrics			
		Best Batch Size	Dice †	mIoU †	FPE ‡	FNE ‡
FCN8	ResNet50 VGG	2	0.2621	0.1914	0.6789	0.6062
		2	0.4993	0.3963	0.4576	<i>0.3824</i>
		2	0.5101	0.3952	0.3643	0.4500
FCN32	ResNet50 VGG	2	0.2174	0.1594	0.7564	0.6980
		2	0.4334	0.3372	0.5090	0.5081
		2	<i>0.5708</i>	<i>0.4549</i>	0.3396	0.3833
SegNet	ResNet50 VGG	32	0.2677	0.1880	0.6325	0.6510
		32	0.4768	0.3676	0.4325	0.4339
		32	0.4596	0.3469	0.4003	0.5158
U-Net	ResNet50 VGG	2	0.4057	0.3035	0.4900	0.5119
		32	0.0646	0.0371	0.5585	0.9584
		2	0.1446	0.0878	<i>0.2501</i>	0.9020
Proposed	Modified VGG	2	<b>0.7447</b>	<b>0.6467</b>	<b>0.1866</b>	<b>0.3056</b>

metrics are:

False Positive Error (FPE) as in Equation (8):

$$FPE = \frac{FP}{FP + TN} \quad (8)$$

and False Negative Error (FNE) as in Equation (9).

$$FNE = \frac{FN}{FN + TP} \quad (9)$$

where  $FP$  is the total number of false positives in the predictions,  $TN$  is the total number of true negatives and  $FN$  is the total number of false negatives.

## V. RESULTS

Table II shows the results when trained on two types of annotation: manual delineations vs refined contours. The results show that the algorithm did not learn as effectively from the human delineation on the boundary (polygonal outlines). The refined contour consistently demonstrated closer agreement with the machine predictions, without relying on the type of ground truth used for training. Therefore, we use image processing refined contour as ground truth for both train set and test set, for the rest of the paper.

As shown in Table III, many of the available techniques give reasonable results in DFU segmentation. Among the baseline methods, the best performing model was FCN32 with a VGG backbone, with the highest Dice score of 0.5708 and 0.4549 for mIoU. A key factor in this task is the ability of the network to handle images without positive DFU cases (True Negatives), thus we use the FPE metric. In such cases the best performing model is also FCN32 VGG, which shows a high understanding of the surround regions. We observe that most methods that use a higher batch size resulted in significant performance degradation. A contributing factor to this is likely to be background noise present in the images

where the environment can vary significantly between images. Lower batch sizes allowed the system to focus on a case by case basis, allowing the network to slowly learn to ignore the background noise and focus on the wounds.

TABLE IV

COMPARISON OF THE DFUC2022 ENTRIES AND OUR PROPOSED METHOD.

Team	Metrics			
	Dice †	mIoU †	FPE ‡	FNE ‡
yllab	<b>0.7287</b>	0.6252	0.2048	0.2341
LkRobotAI Lab	0.7280	<b>0.6276</b>	0.2154	0.2261
agaldran	0.7263	0.6273	0.2262	<b>0.2210</b>
ADAR-LAB	0.7254	0.6245	<b>0.1847</b>	0.2582
seoyoung	0.7220	0.6208	0.1925	0.2584
FHDO	0.7169	0.6130	0.2145	0.2453
GP_2022	0.6986	0.5921	0.2065	0.2778
DGUT-XP	0.6984	0.5945	0.2523	0.2379
IISlab	0.6975	0.5926	0.2163	0.2734
AGH_MVG	0.6725	0.5690	0.2555	0.2830
Ours	<b>0.7447</b>	<b>0.6467</b>	0.1866	0.3056

Table IV highlights the results for DFUC2022. The top 10 scores demonstrate the challenge of DFU segmentation for a wide range of networks. The team yllab achieved the best score in Dice (0.7287) in which the challenge was based. This was closely followed by LKRobotAI Lab, who achieved the highest mIoU (0.6276) showing a high agreement of prediction and ground truth overlap. The 3rd place team, agaldran, achieved the lowest FNE (0.2210), highlighting that they reduced the amount of falsely predicted DFU pixels, whereas the 4th place team, ADAR-LAB, achieved the best FPE score (0.1847). Our method achieves higher Dice (0.7447) and mIoU (0.6467) scores, showing a high degree of agreement between prediction and ground truth. Additionally, we have a slightly higher FPE (0.1866) when compared to the best performing (0.1847). However, one outlier with our method is that we report lowest performance in FNE (0.3056). Our method demonstrates that



Fig. 6. An example of how the inclusion of dilation smoothing improves predictions in the modified network on full images. From left to right: input image, standard FCN32 VGG and modified FCN32 VGG. *Note: For illustration, images were cropped to focus on DFU region.*



Fig. 7. An example of texture similarity and over down-sampling issues in a DFU prediction. From left to right: original image, standard FCN32 VGG, modified FCN32 VGG and ground truth. *Note: For illustration, images were cropped to focus on DFU region.*

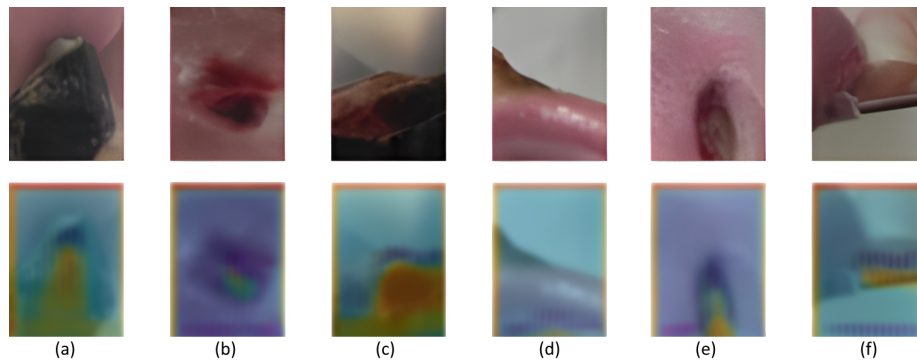


Fig. 8. Examples of the patches used in the modified network, demonstrating the ability of the network to focus on DFU regions, including edge cases (d) and occlusion (f).

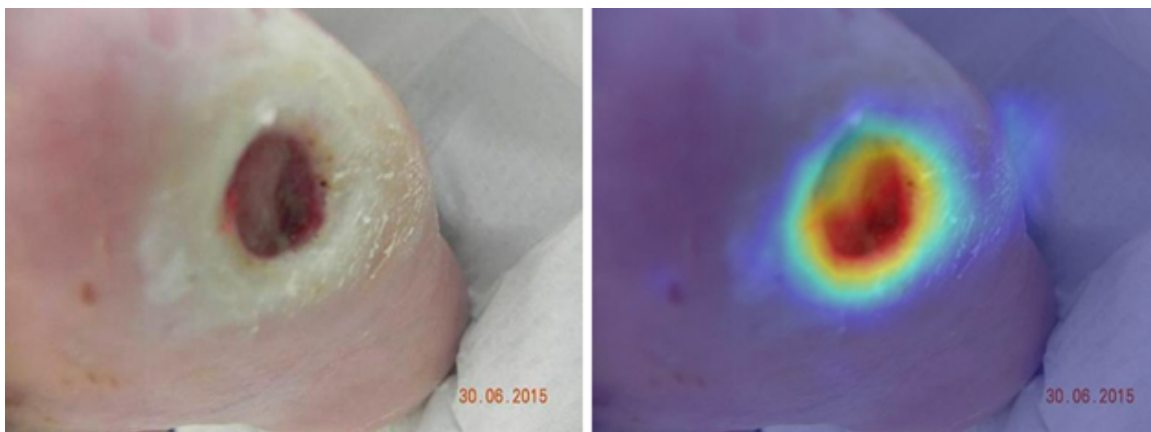


Fig. 9. Illustration of a prediction in a full resolution image, note that the network correctly focuses on DFU regions.

training using patched images and reduced pooling allows for improvements in Dice and mIoU.



The results in Table III and IV show that our proposed strategy and modified network has improved the results and achieved 0.7447 for Dice, 0.6467 for mIoU, 0.1866 for FPE and 0.3056 for FNE. As visually illustrated in Figure 6, the modified network, with the inclusion of dilation smoothing, is able to refine the results within the network. Another example to show the superiority of the modified network is in Figure 7, where due to the similarity between surrounding skin and DFU the standard method fails, but the modified network is able to detect some overlap.

As shown in Figures 8 and 9, the best performing network successfully highlights and focuses on the DFU regions. In addition, these figures highlight how the network modifications allow the system to identify a wide variety of DFU features within an image. However, note that in full size predictions (see Figure 8) the small mark to the left and the damaged skin on the right have also been focused on by the network. This highlights some of the features of the modified network segmentation, similar to Figure 7, the broken skin could indicate a early onset of DFU. Similarly the minor activation on the left could be an indication of a smaller ulcer, due to its colour, shape and texture. Thus, a slight activation over these regions is shown.

## VI. DISCUSSION

We highlight that the best performing baseline methods had several difficulties which reduced model performance, as shown in Figure 10:

- *Excessive down-sampling of images:* Many of the segmentation backbones are based on classification networks in which reducing to core features is essential. However, with the small image to wound size ratio, this removes the full wound from the image.
- *Data distribution:* As show in Figure 2, a large proportion of the dataset has a DFU to background ratio of  $<5\%$ . This represents a large dataset bias towards none-DFU regions. This causes the networks to prioritize on the background class over DFU region, and in some cases the DFU class is ignored.
- *Background noise:* Owing to the shape and location of DFU and patient mobility, many of the images contain a wide assortment of noise. In some cases, the foot is surrounded by a blue or white cloth so the network can focus, but in many cases the background contains clothes, floor details and other medical equipment. This poses a difficulty and the network must learn to cope with a large variety of background data.
- *Region similarity:* With many cases of DFU the textural quality of the lesion is similar to that of the surrounding skin, especially in cases of infection. The textural similarity of DFU regions, periwound and surrounding skin regions, introduces difficulty in distinguishing the regions, as shown in Figure 7. This means that the networks struggled to differentiate between the DFU and other parts of the foot.
- *Rapid up-sampling:* Due to the focus of the backbones ability to output valid feature maps the head of the

network is usually light weight. This results in the up-sampling output being performed at a high rate, causing pixelated regions, in addition to small false detection regions.

These issues are the cause of the difficulties the baseline models produce. Oversampling removes the smaller wounds, which amplifies the problem of data distribution, where most wounds are below 5% of the total image size, meaning the networks focus more on the background than on the DFU regions. Furthermore, this focus on the background data is amplified by the inconstant and noisy data. Owing to this, the region similarity of the DFU, periwound and surrounding skin is made difficult causing some networks to focus on the entire foot over the DFU regions, as there is too much focus on background data. Finally, the networks perform well using the smoothed masks over the original jagged contours provided by clinicians. Thus, in the final stages when re-upsampling to the desired size, pixelation occurs due to the rapid up-sampling, producing block-like segmentation that requires additional post-processing to smooth and remove small regions.

## VII. CONCLUSION

In this paper, we introduce the largest available DFU dataset containing 2000 annotated training images and 2000 test images without annotations, together with the capability of online evaluation of network predictions. We also provide challenging cases, such as non-DFU cases and images resulting from annotator disagreement. We then provide a series of baselines on state-of-the-art models with explainable AI techniques.

We demonstrate that by performing preprocessing on the expert delineation to smooth the DFU regions, the networks were able to produce more accurate DFU segmentation results. This was shown by comparing a cross validation between raw and smoothed masks. From this study we perform an ablation study on widely used semantic segmentation networks, producing a set of baseline results. The prediction results from the trained models highlight the difficulty in DFU clinical delineation where inter reliability can be inconsistent. This work sheds light on the challenges inherent in the development of AI systems which can aid the standardisation of DFU delineation over time to track healing progress.

We identify the shortcomings inherent in traditional segmentation networks and training techniques using the DFUC2022 dataset. From these findings we modified the best performing network and tailor it to the unique challenges presented by the DFU2022 dataset. From these adjustments to the network design we show a significant increase in model performance, without the use of post processing techniques.

Finally, we analyse heatmaps of successfully trained DFU model predictions on DFU regions, which indicate that the network is capable of focusing on ulcer regions and corresponding features when generating final prediction masks. These machine learning advancements will contribute towards supporting healthcare systems to better manage the increasing demands of DFU care, including the accurate and regular monitoring of DFU healing to increase flexibility in treatment plans.

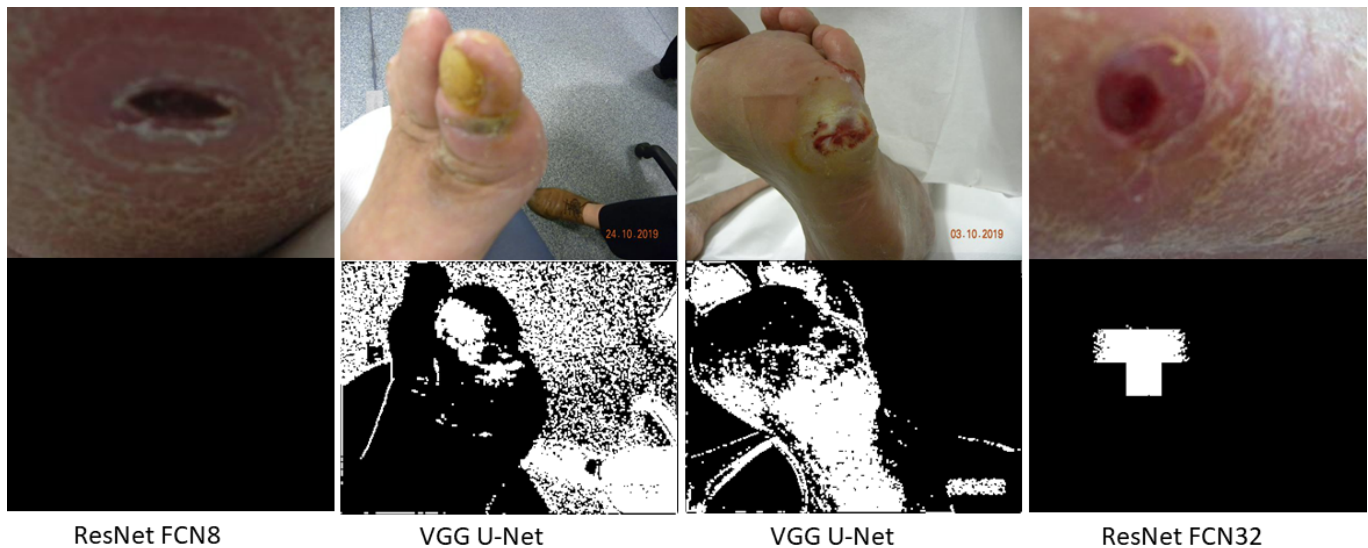


Fig. 10. Illustration of issues associated with the baseline models: (Left) An example over down-sampling removing lesion; (Middle-left) example of background noise effecting prediction; (Middle-right) example of region similarity preventing accurate segmentation; and (Right) example of rapid up-sampling producing a block artefact. *Note: some images were cropped to focus on DFU region.*

## REFERENCES

- [1] E. Ghanassia, L. Villon, J.-F. dit Dieudonne, C. Boegner, A. Avignon, and A. Sultan, "Long-term outcome and disability of diabetic patients hospitalized for diabetic foot ulcers: a 6.5-year follow-up study," *Diabetes care*, vol. 31, no. 7, pp. 1288–1292, 2008.
- [2] W. J. Jeffcoate and K. G. Harding, "Diabetic foot ulcers," *The lancet*, vol. 361, no. 9368, pp. 1545–1551, 2003.
- [3] P. R. Cavanagh, B. A. Lipsky, A. W. Bradbury, and G. Botek, "Treatment for diabetic foot ulcers," *The Lancet*, vol. 366, no. 9498, pp. 1725–1735, 2005.
- [4] K. Ogurtsova, S. Morbach, B. Haastert, M. Dubský, G. Rümenapf, D. Ziegler, A. Jirkovska, and A. Icks, "Cumulative long-term recurrence of diabetic foot ulcers in two cohorts from centres in Germany and the Czech Republic," *Diabetes research and clinical practice*, vol. 172, p. 108621, 2021.
- [5] A. Ahmad, M. Abujbara, H. Jaddou, N. A. Younes, and K. Ajlouni, "Anxiety and depression among adult patients with diabetic foot: prevalence and associated factors," *Journal of clinical medicine research*, vol. 10, no. 5, p. 411, 2018.
- [6] F. M. Davis, A. Kimball, A. Boniakowski, K. Gallagher, and K. Gallagher, "Dysfunctional Wound Healing in Diabetic Foot Ulcers : New Crossroads," 2018.
- [7] M. Edmonds, C. Manu, and P. Vas, "The current burden of diabetic foot disease," *Journal of Clinical Orthopaedics and Trauma*, vol. 17, pp. 88–93, 2021.
- [8] R. Sorber and C. J. Abularrage, "Diabetic foot ulcers: Epidemiology and the role of multidisciplinary care teams," in *Seminars in vascular surgery*, vol. 34, no. 1. Elsevier, 2021, pp. 47–53.
- [9] M. Chang and T. T. Nguyen, "Strategy for treatment of infected diabetic foot ulcers," *Accounts of chemical research*, vol. 54, no. 5, pp. 1080–1093, 2021.
- [10] K. Glover, A. C. Stratakos, A. Varadi, and D. A. Lamprou, "3D scaffolds in the treatment of diabetic foot ulcers: new trends vs conventional approaches," *International Journal of Pharmaceutics*, vol. 599, p. 120423, 2021.
- [11] Z. J. Lo, N. K. Surendra, A. Saxena, and J. Car, "Clinical and economic burden of diabetic foot ulcers: A 5-year longitudinal multi-ethnic cohort study from the tropics," *International Wound Journal*, vol. 18, no. 3, pp. 375–386, 2021.
- [12] H. Sun, P. Saeedi, S. Karuranga, M. Pinkepank, K. Ogurtsova, B. B. Duncan, C. Stein, A. Basit, J. C. N. Chan, J. C. Mbanya, and Others, "IDF diabetes atlas: Global, regional and country-level diabetes prevalence estimates for 2021 and projections for 2045," *Diabetes research and clinical practice*, p. 109119, 2021.
- [13] R. Pranata, J. Henrina, W. M. Raffaello, S. Lawrensia, and I. Huang, "Diabetes and COVID-19: the past, the present, and the future," *Metabolism*, vol. 121, p. 154814, 2021.
- [14] M. H. Yap, R. Hachiuma, A. Alavi, R. Brüngel, B. Cassidy, M. Goyal, H. Zhu, J. Rückert, M. Olshansky, X. Huang, H. Saito, S. Hassanpour, C. M. Friedrich, D. B. Ascher, A. Song, H. Kajita, D. Gillespie, N. D. Reeves, J. M. Pappachan, C. O'Shea, and E. Frank, "Deep learning in diabetic foot ulcers detection: A comprehensive evaluation," *Computers in Biology and Medicine*, vol. 135, p. 104596, 2021. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0010482521003905>
- [15] M. H. Yap, B. Cassidy, J. M. Pappachan, C. O'Shea, D. Gillespie, and N. D. Reeves, "Analysis towards classification of infection and ischaemia of diabetic foot ulcers," in *2021 IEEE EMBS International Conference on Biomedical and Health Informatics (BHI)*, 2021, pp. 1–4.
- [16] M. H. Yap, R. Hachiuma, A. Alavi, R. Brüngel, B. Cassidy, M. Goyal, H. Zhu, J. Rückert, M. Olshansky, X. Huang, H. Saito, S. Hassanpour, C. M. Friedrich, D. B. Ascher, A. Song, H. Kajita, D. Gillespie, N. D. Reeves, J. M. Pappachan, C. O'Shea, and E. Frank, "Deep learning in diabetic foot ulcers detection: A comprehensive evaluation," *Computers in Biology and Medicine*, vol. 135, p. 104596, 2021. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0010482521003905>
- [17] B. Cassidy, C. Kendrick, N. Reeves, J. Pappachan, C. O'Shea, D. Armstrong, and M. H. Yap, *Diabetic Foot Ulcer Grand Challenge 2021: Evaluation and Summary*, 01 2022, pp. 90–105.
- [18] M. H. Yap, C. Kendrick, N. Reeves, M. Goyal, J. Pappachan, and B. Cassidy, *Development of Diabetic Foot Ulcer Datasets: An Overview*, 01 2022, pp. 1–18.
- [19] N. D. Reeves, B. Cassidy, C. A. Abbott, and M. H. Yap, "Chapter 7 - novel technologies for detection and prevention of diabetic foot ulcers," in *The Science, Etiology and Mechanobiology of Diabetes and its Complications*, A. Gefen, Ed. Academic Press, 2021, pp. 107–122. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/B9780128210703000076>
- [20] B. Cassidy, N. D. Reeves, J. M. Pappachan, N. Ahmad, S. Haycocks, D. Gillespie, and M. H. Yap, "A cloud-based deep learning framework for remote detection of diabetic foot ulcers," *arXiv preprint arXiv:2004.11853*, 2021.
- [21] C. Wang, D. M. Anisuzzaman, V. Williamson, M. K. Dhar, B. Rostami, J. Niezgoda, S. Gopalakrishnan, and Z. Yu, "Fully automatic wound segmentation with deep convolutional neural networks," *Scientific Reports*, vol. 10, no. 1, pp. 1–9, 2020. [Online]. Available: <https://doi.org/10.1038/s41598-020-78799-w>
- [22] S. Thomas, "Medetec," 2020, last access: 08/11/21. [Online]. Available: <http://www.medetec.co.uk/index.html>
- [23] C. Wang, B. Rostami, J. Niezgoda, S. Gopalakrishnan, and Z. Yu, "Foot ulcer segmentation challenge 2021," Mar. 2021. [Online]. Available: <https://doi.org/10.5281/zenodo.4575314>
- [24] C. Wang, A. Mahbod, I. Ellinger, A. Galdran, S. Gopalakrishnan, J. Niezgoda, and Z. Yu, "Fuseg: The foot ulcer segmentation challenge," 2022. [Online]. Available: <https://arxiv.org/abs/2201.00414>

- [25] 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)*, Oct 2017, pp. 618–623.
- [26] A. Mahbod, R. Ecker, and I. Ellinger, "Automatic foot ulcer segmentation using an ensemble of convolutional neural networks," 2021. [Online]. Available: <https://arxiv.org/abs/2109.01408>
- [27] D.-J. Kroon, "Snake: Active contour," Online, Jan. 2022. [Online]. Available: <https://www.mathworks.com/matlabcentral/fileexchange/28149-snake-active-contour>
- [28] M. Goyal and M. H. Yap, "Multi-class semantic segmentation of skin lesions via fully convolutional networks," *arXiv preprint arXiv:1711.10449*, 2017.
- [29] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, "Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs," *IEEE transactions on pattern analysis and machine intelligence*, vol. 40, no. 4, pp. 834–848, 2017.
- [30] Z. Khan, N. Yahya, K. Alsaih, S. S. A. Ali, and F. Meriaudeau, "Evaluation of deep neural networks for semantic segmentation of prostate in t2w mri," *Sensors*, vol. 20, no. 11, 2020. [Online]. Available: <https://www.mdpi.com/1424-8220/20/11/3183>
- [31] R. Azad, M. Asadi-Aghbolaghi, M. Fathy, and S. Escalera, "Attention deeplabv3+: Multi-level context attention mechanism for skin lesion segmentation," in *Computer Vision – ECCV 2020 Workshops*, A. Bartoli and A. Fusiello, Eds. Cham: Springer International Publishing, 2020, pp. 251–266.
- [32] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 3431–3440.
- [33] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2015, pp. 234–241.
- [34] V. Badrinarayanan, A. Kendall, and R. Cipolla, "Segnet: A deep convolutional encoder-decoder architecture for image segmentation," *IEEE transactions on pattern analysis and machine intelligence*, vol. 39, no. 12, pp. 2481–2495, 2017.
- [35] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [36] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [37] M. D. Zeiler, "ADADELTA: An Adaptive Learning Rate Method," 2012. [Online]. Available: <http://arxiv.org/abs/1212.5701>
- [38] Y. N. Dauphin, A. Fan, M. Auli, and D. Grangier, "Language modeling with gated convolutional networks," in *International conference on machine learning*. PMLR, 2017, pp. 933–941.
- [39] J. Hu, L. Shen, and G. Sun, "Squeeze-and-excitation networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 7132–7141.