

# Quantifying the Effect of Image Similarity on Diabetic Foot Ulcer Classification

Imran Chowdhury Dipto<sup>1</sup>[0000-0001-9183-2651], Bill Cassidy<sup>1\*</sup>[0000-0003-3741-8120], Connah Kendrick<sup>1</sup>[0000-0002-3623-6598], Neil D. Reeves<sup>3</sup>[0000-0001-9213-4580], Joseph M. Pappachan<sup>2</sup>[0000-0003-0886-5255], Vishnu Chandrabalan<sup>2</sup>[0000-0002-2687-1096], and Moi Hoon Yap<sup>1</sup>[0000-0001-7681-4287]

<sup>1</sup> Centre for Advanced Computational Science, Department of Computing and Mathematics, Manchester Metropolitan University, Manchester M1 5GD, United Kingdom

<sup>2</sup> Lancashire Teaching Hospitals NHS Trust, Preston, PR2 9HT, United Kingdom

<sup>3</sup> Musculoskeletal Science and Sports Medicine Research Centre, Manchester Metropolitan University, Manchester M1 5GD, United Kingdom  
M.Yap@mmu.ac.uk

**Abstract.** This research conducts an investigation on the effect of visually similar images within a publicly available diabetic foot ulcer dataset when training deep learning classification networks. The presence of binary-identical duplicate images in datasets used to train deep learning algorithms is a well known issue that can introduce unwanted bias which can degrade network performance. However, the effect of visually similar non-identical images is an under-researched topic, and has so far not been investigated in any diabetic foot ulcer studies. We use an open-source fuzzy algorithm to identify groups of increasingly similar images in the Diabetic Foot Ulcers Challenge 2021 (DFUC2021) training dataset. Based on each similarity threshold, we create new training sets that we use to train a range of deep learning multi-class classifiers. We then evaluate the performance of the best performing model on the DFUC2021 test set. Our findings show that the model trained on the training set with the 80% similarity threshold images removed achieved the best performance using the InceptionResNetV2 network. This model showed improvements in F1-score, precision, and recall of 0.023, 0.029, and 0.013, respectively. These results indicate that highly similar images can contribute towards the presence of performance degrading bias within the Diabetic Foot Ulcers Challenge 2021 dataset, and that the removal of images that are 80% similar from the training set can help to boost classification performance.

## 1 Introduction

Since the publication of the DFUC2021 Proceedings, there has been no substantial progress made on the DFU multi-class classification task. This paper

---

\* equal contribution

studies one of possible cause, i.e., the effect of visually similar non-identical images within DFUC2021 dataset. Image duplication (the presence of binary identical images) is generally acknowledged as a factor in reducing model performance when training deep learning models due to the performance degrading bias that over-represented features may introduce into the trained model. However, the effect of visually similar non-identical images which may result in an over-representation of certain features present in deep learning datasets is an under-researched topic. An overabundance of certain features in a dataset may cause undesirable performance degrading bias in any models trained using them. In this paper we conduct an analysis of the effect of images that are visually similar but not binary identical on a publicly available diabetic foot ulcer dataset using an open-source fuzzy matching algorithm. We train a large range of multi-class deep learning classification models on the Diabetic Foot Ulcer Challenge 2021 dataset (DFUC2021) [1], and for the first time, quantify the effect of image similarity on network accuracy.

We found no studies that observed and quantified the effect of feature over-representation or the effect of image similarity in DFU research. The effect of binary duplicate images has been observed in other domains [2] but the topic remains an under-researched problem generally. A common theme with many previous studies is limited dataset size. A small dataset may hinder the ability of models to generalise to a wider range of examples in real-world settings. Conversely, a large dataset might introduce performance degrading feature-bias if the data has been collected from a small number of subjects. To address this, we conduct experiments to analyse the effect of image similarity on the DFUC2021 dataset.

## 2 Related Work

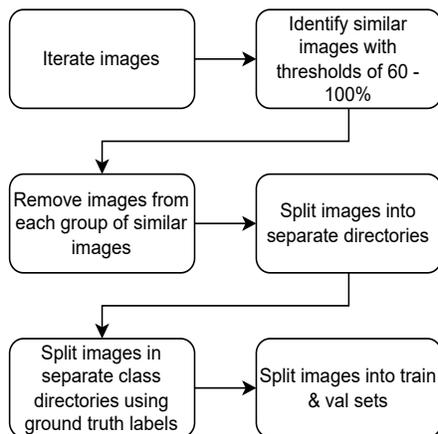
Previous research on DFU has involved localisation [3,4,5,6,7,8], binary [9] and multi-class classification [10,11], and segmentation [12,13]. Goyal et al. [14] proposed a deep learning architecture for DFU classification. Their work was notable for achieving high scores in sensitivity and accuracy using a small DFU dataset (< 2000 images for training and testing).

More recently, Al-Garaawi et al. [15] conducted a series of binary classification experiments using DFU patches. This work used mapped local binary pattern coded images with RGB images as inputs to the CNN to increase binary classifier performance. However, a limitation of this work is the use of small datasets.

In last year’s Diabetic Foot Ulcers Grand Challenge, Yap et al. [1] conducted multi-class classification experiments using the DFUC2021 dataset. This work highlighted the challenging nature of multi-class classification in this domain due to intra-class similarities.

### 3 Methodology

To observe the effect of similar images in the DFUC2021 dataset on multi-class classification, we devised a strategy of gradually removing successive groups of similar images from the training set. Each group of similar images were identified using the dupeGuru [16] Windows application. This open-source application implements a fuzzy search algorithm capable of identifying visually similar images. Results can be filtered by percentage similarity within the application. Using this feature, we were able to identify groups of similar images within the DFUC2021 training set. For each similarity threshold, we train a set of multi-class classifiers capable of classifying the following five classes: (1) control, (2) infection, (3) ischemia, (4) infection and ischemia, and (5) none. Figure 1 shows an overview of the entire process used to create the new training sets used in our experiments.



**Fig. 1.** Overview of the process used in the identification and removal of similar images at each similarity threshold.

#### 3.1 Dataset Description

For our experiments, we use the publicly available DFUC2021 dataset, introduced by Yap et al. [1]. This dataset is the largest publicly available DFU dataset with wound pathology class labels. The dataset comprises a total of 15,683 images, sized at  $224 \times 224$  pixels, with 5955 images for the training set (2555 infection only, 227 ischaemia only, 621 both infection and ischaemia, and 2552 without ischaemia and infection), 3994 unlabeled images, and 5734 images for the testing set. All wounds are cropped from larger images so that only the wound is present in each image. The DFUC2021 dataset is highly heterogeneous due to the nature of the variety of capture devices and variable settings used

during photographic acquisition. The ground truth labels were provided by expert clinicians at Lancashire Teaching Hospitals, UK. This dataset was obtained with ethical approval from the UK National Health Service Research Ethics Committee (reference number: 15/NW/0539).

### 3.2 Fuzzy Algorithm

The fuzzy algorithm used by the dupeGuru application reads each image in RGB bitmap mode which is then split into blocks. Next, the analysis phase uses a  $15 \times 15$  pixel grid to average the colour of each grid tile, the results of which are stored in a cache database. Each grid tile, representing an average colour, is then compared to its corresponding grid on the other image being compared to, and a sum of the difference between R, G and B on each side is computed. The RGB sums are then added together to obtain a final result. If the score is smaller or equal to the user-specified threshold, then a match is found. If a threshold of 100 is set by the user then the algorithm adds an extra constraint indicating that images should contain identical binary data.

**Table 1.** Summary of the number of similar images found by the dupeGuru fuzzy algorithm in the train, test, and train & test sets combined at each user-defined similarity threshold.

Threshold (%)	Similar Images		
	Train Set	Test Set	Train & Test Set
100	0	0	0
95	1	0	3
90	19	23	48
85	106	125	268
80	317	345	719
75	590	621	1278
70	1013	979	2082
65	1509	1367	2976
60	2066	1683	3906

### 3.3 Identification of Similar Images

To identify the similar images in the DFUC2021 dataset we have used the hardness level (similarity threshold) in the dupeGuru application with the values of 60%, 65%, 70%, 75%, 80%, 85%, 90%, 95%, and 100%. A similarity threshold of 80%, for example, indicates that the application will find images that have 80% similarity. We ran the fuzzy algorithm on the training, test, and the training and test sets combined to find similar images that exist exclusively within the training set, exclusively within the test set, and within both train and test sets combined. Table 1 shows a summary of the number of similar images detected

by the dupeGuru fuzzy algorithm on the different thresholds in the training set, the test set, and the training and test sets combined.

### 3.4 Removal of Similar Images as Determined by Similarity Thresholds

Images of each of the classes present in the dataset were separated into different directories - one directory per class. On each of these directories the dupeGuru fuzzy algorithm was run using similarity thresholds of 60% to 100%. Next, the filenames for the images in each similarity threshold were saved into CSV files containing Group ID and Image filenames. Group ID refers to the grouping of similar images returned by the fuzzy algorithm, where each group of similar images is assigned a unique sequential identifier. The CSV files output by dupeGuru were then merged with the file containing the ground truth labels of the training set based on the image filenames. For each group of similar images, all but the first image in each group was removed. This meant that a single example from each similarity group was kept for inclusion in each similarity threshold training set.

**Table 2.** Summary of the number of similar images removed at each similarity threshold and the remaining images that are used for the new training sets.

Threshold (%)	Similar Images Removed	Remaining Images
95	1	9948
90	19	9930
85	106	9843
80	317	9632
75	590	9359
70	1013	8936
65	1508	8441
60	2068	7881

By comparing the filenames with the ground truth labels, the images in the curated datasets were copied into new directories which formed the new training sets. To check the validity of the results from this process, a Python routine was created which compared the CSV files against the ground truth labels to ensure that the correct images had been copied and that non of the additional images from each similarity group were present in the new training sets. An additional manual spot-check was completed on a random sample of images in the new training sets to ensure that the images had been correctly separated. Table 2 shows a summary of the number of images removed at each similarity threshold together with the total remaining images used to form the new training sets.

## 4 Image Similarity Analysis

In this section we analyse a selection of images from each of the similarity thresholds returned by the dupeGuru fuzzy algorithm prior to training the multi-class classification models. Note that not all similarity searches returned results.

### 4.1 Train Set Image Similarity

Figure 2 (a & b) shows two images from the training set in the 75% similarity threshold. These images are of the same wound at different levels of magnification, with the example shown in (a) being at a higher level of magnification. Figure 2 (c & d) shows two training set images in the 65% similarity threshold. These two images represent two distinctly different DFU wounds with noticeably different features. Figure 2 (e & f) shows a further two training set images which were identified in the 60% similarity threshold. As with the previous examples, these images represent two different wounds.

### 4.2 Test Set Image Similarity

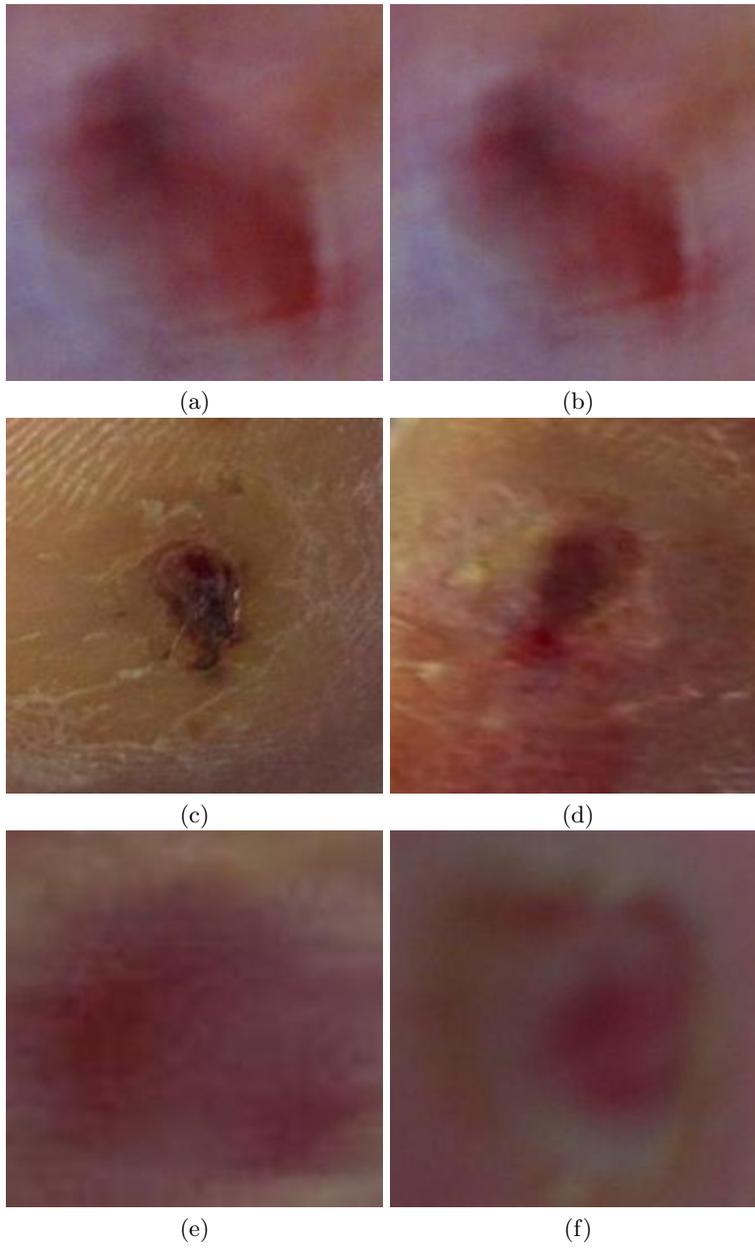
Figure 3 (a & b) shows two similar images from the 80% threshold. These images are of the same wound, with the second image being a natural augmentation case with a slightly different zoom level. The main visual differences can be observed on the bottom section of the image where the dark spots in image (a) are not present on image (b). Figure 3 (c & d) shows two similar images from the test set with 65% similarity. As per the previous test examples (Figure 3 (a & b)), the second image represents a case of natural augmentation where the wound has a noticeably increased zoom level.

### 4.3 Train & Test Set Image Similarity

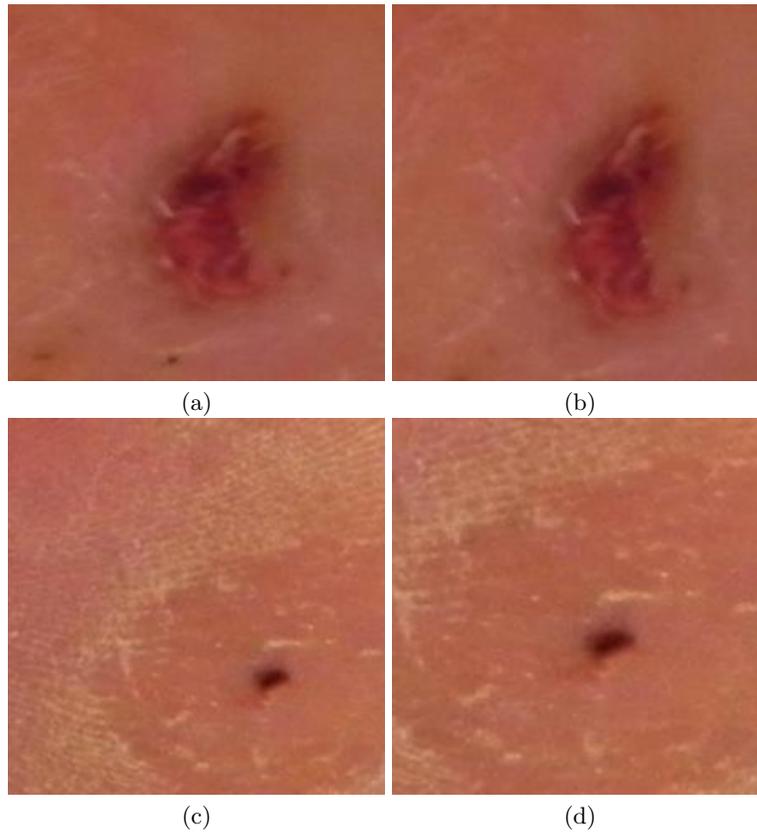
Figure 4 (a & b) shows two distinctly different wound images identified in the 75% similarity threshold. Image (a) is from the training set, image (b) is from the test set. Figure 4 (c & d) and (e & f) show further examples of visually similar images found across training and test sets at 65% and 60% similarity thresholds respectively. Note that we do not discuss the class of the test set examples as these are part of a live public challenge for DFUC2021 which is still open to submissions (<https://dfu-challenge.github.io/dfuc2021.html>).

### 4.4 Inter-class Image Similarity

Our experiments using the dupeGuru fuzzy algorithm did not return any inter-class similarity results for the following groups of classes: (1) ‘both’ vs ‘none’, and (2) ‘infection’ vs ‘ischemia’. For the ‘infection’ vs ‘none’ similarity searches, similar images were found for the 70% (45 images), 75% (13 images), and 80% (4 images) similarity brackets. For the ‘ischemia’ vs ‘none’ similarity searches,

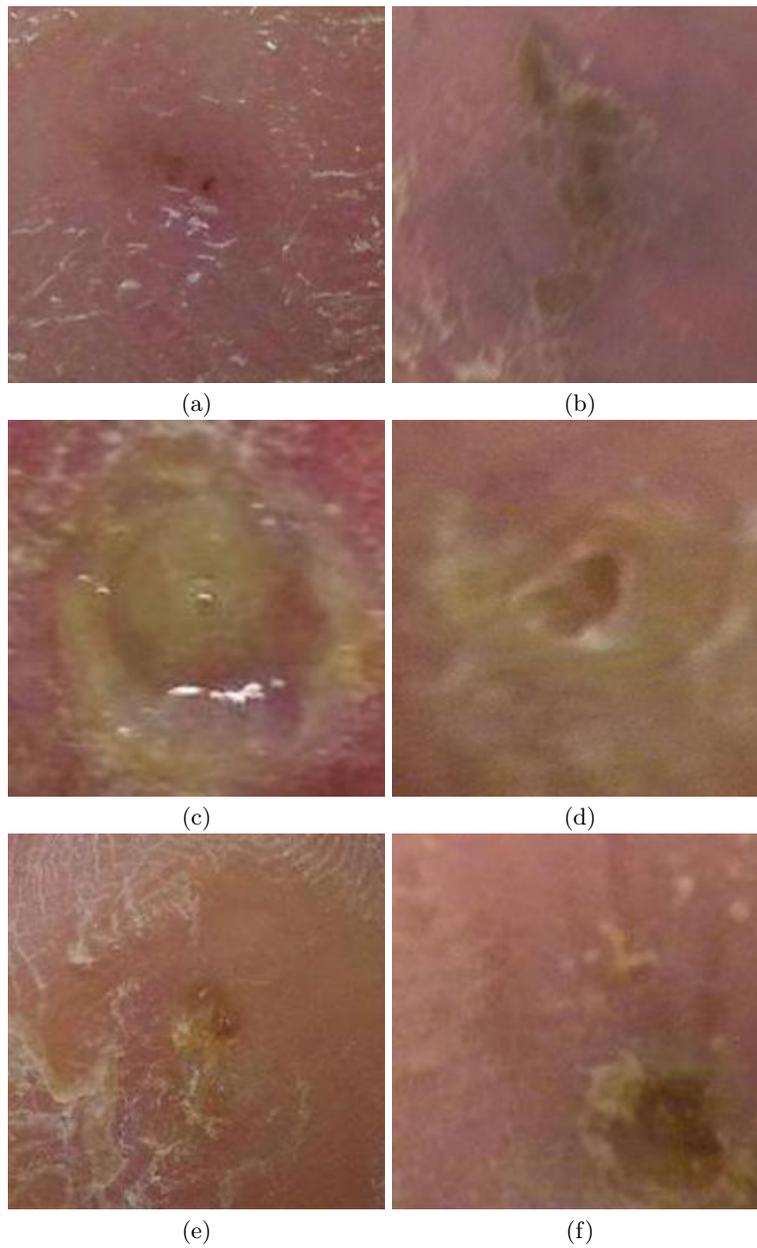


**Fig. 2.** Illustration of two training set images identified by the dupeGuru fuzzy algorithm with a similarity threshold of 75% (a & b), two training set images from the 'none' class found in the 65% similarity threshold (c & d), and two training set images found in the 60% similarity threshold - (e) is from the 'none' class, (f) is from the 'unlabelled' class.

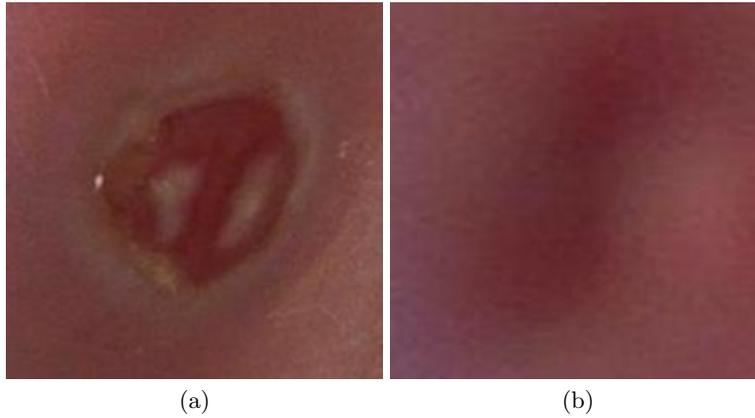


**Fig. 3.** Illustration of two test set images identified by the dupeGuru fuzzy algorithm with a similarity threshold of 80% (a & b), and two test set images identified with a similarity threshold of 65% (c & d).

similar images were found for the 70% (5 images) and 75% (3 images) brackets. Figure 5 shows two images from the 75% similarity threshold training set, with image (a) showing an ‘infection’ class example, and image (b) showing a ‘none’ class example. Figure 6 shows two images from the 70% similarity threshold training set, where image (a) shows an ‘ischemia’ class example, while image (b) shows a ‘none’ class example.



**Fig. 4.** Illustration of similar images located across both training and test sets identified in the 60% to 75% similarity thresholds. Image (a) shows a training set image and image (b) shows a test set image, both identified in the 75% similarity threshold. Image (c) shows a training set image and image (d) shows a test set image, both identified in the 65% similarity threshold. Image (e) shows a training set image and image (f) shows a test set image, both identified in the 60% similarity threshold.



**Fig. 5.** Illustration of two images from the training set which were identified in the inter-class similarity results for the 75% similarity threshold: (a) an image from the ‘infection’ class, and (b) an image from the ‘none’ class



**Fig. 6.** Illustration of inter-class similarity results in the 70% similarity threshold: (a) an image from the ‘ischemia’ class, and (b) an image from the ‘none’ class.

#### 4.5 Model Training

Following the creation of each training set, as per Table 2, we trained a selection of popular deep learning classification networks using each of our curated training sets. Batch size was set to 32 with stochastic gradient descent used for the optimiser and categorical cross-entropy as the loss function. Early stopping was implemented monitoring validation accuracy with a patience of 10. The hardware configuration used for our experiments was as follows: Intel Core i7-10750H CPU @2.60GHz, 64GB RAM, NVIDIA GeForce RTX 2070 Super with Max-Q Design

8GB. The software configuration used was as follows: Ubuntu 20.04 LTS, Python 3.8.13, and Tensorflow 2.4.2.

## 5 Results and Discussion

This section details the results of training the multi-class classification networks on each of the training sets as detailed in Table 2. We report the validation results from the full training set and the training sets using the following similarity thresholds: (1) 60%, (2) 65%, (3) 70%, (4) 75%, (5) 80%, (6) 85%, (7) 90%, and (8) 95%. Finally, we report the test results for the model with the highest validation accuracy.

### 5.1 Baseline Results

Table 3 shows the validation accuracy results for the models trained using the full training set, with no similar images removed. The InceptionResNetV2 model shows a clear lead in validation accuracy with 0.801, showing an increase of 0.038 over the next best performing model, which was ResNet50 with a validation accuracy of 0.763.

**Table 3.** Best epoch and validation accuracy of the models trained on the full DFUC2021 dataset.

Model	Best Epoch	Validation Accuracy
DenseNet201	36	0.731
EfficientNetB0	25	0.656
EfficientNetB1	11	0.587
EfficientNetB3	38	0.649
<b>InceptionResNetV2</b>	42	<b>0.801</b>
InceptionV3	16	0.704
<b>ResNet50</b>	47	<b>0.752</b>
ResNet50V2	43	0.738
ResNet101	30	0.719
ResNet101V2	54	0.735
ResNet152	19	0.703
<b>ResNet152V2</b>	76	<b>0.763</b>
VGG16	52	0.687
VGG19	75	0.716
Xception	18	0.735

### 5.2 Results on the Curated Datasets

The validation results for each of our curated training sets are shown in Table 4, Table 5, Table 6, Table 7, Table 8, Table 9, Table 10, and Table 11. The

best performing model for validation accuracy is InceptionResNetV2 on the 80% similarity threshold (0.885), as shown in Table 8. This represents an increase of 0.030 over the next best performing model across all the best performing models in all similarity thresholds, which was InceptionResNetV2 in the 75% similarity threshold with 0.855 accuracy as shown in Table 7.

**Table 4.** Best Epoch and validation accuracy of the models trained on the 60% similarity threshold dataset.

Model	Best Epoch	Validation Accuracy
DenseNet121	2	0.399
DenseNet169	2	0.419
DenseNet201	3	0.403
EfficientNetB0	5	0.445
EfficientNetB1	4	0.424
EfficientNetB3	7	0.443
<b>InceptionResNetV2</b>	3	<b>0.470</b>
<b>InceptionV3</b>	3	<b>0.461</b>
ResNet50	3	0.431
ResNet50V2	8	0.416
ResNet101	9	0.410
ResNet101V2	1	0.421
<b>ResNet152</b>	2	<b>0.460</b>
ResNet152V2	2	0.432
VGG16	10	0.435
VGG19	14	0.442
Xception	5	0.416

The lowest performing models across all similarity thresholds for validation accuracy were ResNet152 (0.460), InceptionV3 (0.461), and InceptionResNetV2 (0.470), all present in the 60% similarity threshold (see Table 4). This indicates that the 60% similarity threshold removed too many useful examples that the models were able to learn from - 2066 images compared to 317 images for the best performing network (InceptionResNetV2 at 80% similarity threshold) in validation accuracy. All models trained in the 60% similarity threshold also show low convergence for best epoch when compared to all other similarity thresholds, further highlighting a lack of learnable features present in this heavily curated training set.

Given that the InceptionResNetV2 model trained on the 80% similarity threshold training set performed best in validation accuracy, we used this model to obtain test results on the DFUC2021 testing set. The results for this experiment are presented in Table 12. The InceptionResNetV2 model trained on the 80% similarity threshold training set shows clear performance improvements for macro average F1-score, precision, and recall, with improvements of 0.023, 0.029, 0.013 respectively. The reported AUC is slightly higher for the model trained on

**Table 5.** Best Epoch and validation accuracy of the models trained on the 65% similarity threshold dataset.

Model	Best Epoch	Validation Accuracy
DenseNet121	19	0.673
DenseNet169	33	0.731
DenseNet201	21	0.715
EfficientNetB0	44	0.670
EfficientNetB1	53	0.664
EfficientNetB3	28	0.662
<b>InceptionResNetV2</b>	54	<b>0.823</b>
<b>InceptionV3</b>	37	<b>0.756</b>
<b>ResNet50</b>	54	<b>0.754</b>
ResNet50V2	35	0.742
ResNet101	32	0.729
ResNet101V2	17	0.671
ResNet152	26	0.713
ResNet152V2	55	0.744
VGG16	21	0.651
VGG19	42	0.662
Xception	26	0.733

**Table 6.** Best Epoch and validation accuracy of the models trained on the 70% similarity threshold dataset.

Model	Best Epoch	Validation Accuracy
<b>DenseNet121</b>	60	<b>0.765</b>
DenseNet169	34	0.707
DenseNet201	15	0.686
EfficientNetB0	30	0.693
EfficientNetB1	27	0.660
EfficientNetB3	7	0.581
<b>InceptionResNetV2</b>	29	<b>0.759</b>
InceptionV3	26	0.715
ResNet50	7	0.650
ResNet50V2	18	0.684
ResNet101	27	0.719
ResNet101V2	28	0.684
ResNet152	27	0.707
ResNet152V2	14	0.660
VGG16	28	0.669
VGG19	34	0.686
<b>Xception</b>	48	<b>0.791</b>

**Table 7.** Best Epoch and validation accuracy of the models trained on the 75% similarity threshold dataset.

Model	Best Epoch	Validation Accuracy
DenseNet121	27	0.734
DenseNet169	12	0.692
DenseNet201	16	0.684
EfficientNetB0	56	0.704
EfficientNetB1	17	0.634
EfficientNetB3	71	0.73
<b>InceptionResNetV2</b>	93	<b>0.855</b>
InceptionV3	18	0.717
ResNet50	24	0.706
<b>ResNet50V2</b>	45	<b>0.747</b>
<b>ResNet101</b>	32	<b>0.74</b>
ResNet101V2	19	0.675
ResNet152	28	0.706
ResNet152V2	61	0.743
VGG16	50	0.70
VGG19	58	0.712
Xception	17	0.734

**Table 8.** Best epoch and validation accuracy of the models trained on the 80% similarity threshold dataset.

Model	Best Epoch	Validation Accuracy
DenseNet121	8	0.662
DenseNet169	44	0.759
DenseNet201	20	0.704
EfficientNetB0	60	0.708
EfficientNetB1	15	0.619
EfficientNetB3	34	0.681
<b>InceptionResNetV2</b>	99	<b>0.885</b>
<b>InceptionV3</b>	43	<b>0.769</b>
ResNet50	43	0.765
ResNet50V2	32	0.722
ResNet101	46	0.747
<b>ResNet101V2</b>	86	<b>0.789</b>
ResNet152	45	0.733
ResNet152V2	38	0.736
VGG16	25	0.674
VGG19	22	0.672
Xception	30	0.765

**Table 9.** Best epoch and validation accuracy of the models trained on the 85% similarity threshold dataset.

Model	Best Epoch	Validation Accuracy
DenseNet121	15	0.688
DenseNet169	22	0.709
DenseNet201	27	0.742
EfficientNetB0	47	0.700
EfficientNetB1	14	0.589
EfficientNetB3	16	0.605
<b>InceptionResNetV2</b>	52	<b>0.805</b>
InceptionV3	28	0.707
ResNet50	53	0.748
<b>ResNet50V2</b>	64	<b>0.767</b>
ResNet101	35	0.717
ResNet101V2	50	0.735
<b>ResNet152</b>	61	<b>0.755</b>
ResNet152V2	22	0.683
VGG16	61	0.683
VGG19	44	0.682
Xception	13	0.715

**Table 10.** Best epoch and validation accuracy of the models trained on the 90% similarity threshold dataset.

Model	Best Epoch	Validation Accuracy
DenseNet121	20	0.693
DenseNet169	36	0.740
DenseNet201	14	0.700
EfficientNetB0	20	0.653
EfficientNetB1	40	0.658
EfficientNetB3	40	0.659
<b>InceptionResNetV2</b>	41	<b>0.785</b>
<b>InceptionV3</b>	57	<b>0.768</b>
ResNet50	35	0.736
ResNet50V2	23	0.699
<b>ResNet101</b>	46	<b>0.743</b>
ResNet101V2	10	0.665
ResNet152	28	0.715
ResNet152V2	10	0.649
VGG16	61	0.698
VGG19	61	0.689
Xception	6	0.669

**Table 11.** Best epoch and validation accuracy of the models trained on the 95% similarity threshold dataset.

Model	Best Epoch	Validation Accuracy
<b>DenseNet121</b>	78	<b>0.817</b>
<b>DenseNet169</b>	34	<b>0.740</b>
DenseNet201	24	0.720
EfficientNetB0	34	0.677
EfficientNetB1	31	0.664
EfficientNetB3	12	0.595
InceptionResNetV2	21	0.734
InceptionV3	28	0.726
ResNet50	46	0.737
ResNet50V2	39	0.726
ResNet101	38	0.728
ResNet101V2	51	0.737
<b>ResNet152</b>	107	<b>0.820</b>
ResNet152V2	15	0.690
VGG16	37	0.677
VGG19	70	0.715
Xception	22	0.730

the full training set, however, this value is negligible with a difference of just 0.001.

**Table 12.** Macro average test performance metrics for the InceptionResNetV2 model trained on the full DFUC2021 training set and the InceptionResNetV2 model trained on the 80% similarity threshold training set. AUC - area under the curve.

Model	F1-Score	Precision	Recall	AUC
InceptionResNetV2 (full)	0.511	0.523	0.541	<b>0.841</b>
InceptionResNetV2 (80)	<b>0.534</b>	<b>0.552</b>	<b>0.554</b>	0.840

The F1-scores for the multi-class test performance of the InceptionResNetV2 model trained on the 80% similarity threshold training set are shown in Table 13. The InceptionResNetV2 model trained on the 80% similarity threshold training set shows a clear performance increase for all classes, with improvements of 0.011 for the none class, 0.025 for the infection class, 0.015 for the ischemia class, and 0.039 for the both class (infection and ischemia). The biggest performance increase is shown for the both class (0.039). Accuracy for all classes is 0.602 for the model trained on the full DFUC2021 training set, and 0.621 for the model trained on the 80% similarity threshold training set. This demonstrates an accuracy improvement of 0.019 when testing using the InceptionResNetV2 model trained on the 80% similarity threshold training set.

**Table 13.** F1-score multi-class test results for the InceptionResNetV2 model trained on the full DFUC2021 training set and the InceptionResNetV2 model trained on the 80% similarity threshold training set.

Model	None	Infection	Ischemia	Both	Accuracy
Full	0.707	0.512	0.431	0.394	0.602
80%	<b>0.718</b>	<b>0.537</b>	<b>0.446</b>	<b>0.433</b>	<b>0.621</b>

We observe that a number of the visually similar images identified by the dupeGuru fuzzy algorithm were examples of natural augmentation. Our findings indicate that the excess use of subtle augmentation cases does not have the desired effect of boosting network performance. This highlights the importance of rigorously experimenting using individual augmentation sets when training deep learning networks to ascertain if models are being negatively affected by certain augmentation types. We encourage researchers working in other deep learning domains to follow these guidelines in future work to ensure that models are effectively trained, and that the effect of individual augmentation types is better understood.

Our experiments focused on the use of a single image similarity algorithm - an open-source fuzzy algorithm found in the dupeGuru application. Future research might test other image similarity methods, such as the structural similarity index measure, cosine similarity, or mean squared error [2].

## 6 Conclusion

In this work we observed and quantified the effect of non-identical similar images on a selection of popular deep learning multi-class classification networks trained using a large publicly available diabetic foot ulcer dataset. We found that model accuracy is negatively affected by the presence of non-identical visually similar images, but that the removal of too many non-identical visually similar images can degrade network performance. We report our findings to encourage researchers to experiment with other deep learning datasets to gauge a better understanding of the effect of image similarity and the potential bias it may introduce into models trained on such data.

## Acknowledgment

We gratefully acknowledge the support of NVIDIA Corporation who provided access to GPU resources for the DFUC2020 and DFUC2021 Challenges.

## References

1. Moi Hoon Yap, Bill Cassidy, Joseph M Pappachan, Claire O’Shea, David Gillespie, and Neil 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)*, pages 1–4. IEEE, 2021.
2. 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*, 75:102305, 2022.
  3. Manu Goyal, Neil Reeves, Satyan Rajbhandari, and Moi Hoon Yap. Robust methods for real-time diabetic foot ulcer detection and localization on mobile devices. *IEEE journal of biomedical and health informatics*, 2018.
  4. 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.
  5. 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, An-ping 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.
  6. 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.
  7. 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. *IEEE Pervasive Computing*, (01):1–9, jan 2022.
  8. Joseph M Pappachan, Bill Cassidy, Cornelius James Fernandez, Vishnu Chandrabalan, and Moi Hoon Yap. The role of artificial intelligence technology in the care of diabetic foot ulcers: the past, the present, and the future. *World Journal of Diabetes*, 13:1131–1139, 12 2022.
  9. 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.
  10. Bill Cassidy, Connah Kendrick, Neil Reeves, Joseph Pappachan, Claire O’Shea, David Armstrong, and Moi Hoon Yap. *Diabetic Foot Ulcer Grand Challenge 2021: Evaluation and Summary*, pages 90–105. 01 2022.
  11. Moi Hoon Yap, Connah Kendrick, Neil Reeves, Manu Goyal, Joseph Pappachan, and Bill Cassidy. *Development of Diabetic Foot Ulcer Datasets: An Overview*, pages 1–18. 01 2022.
  12. Manu Goyal, Moi Hoon Yap, Neil D Reeves, Satyan Rajbhandari, and Jennifer Spragg. Fully convolutional networks for diabetic foot ulcer segmentation. In *2017 IEEE international conference on systems, man, and cybernetics (SMC)*, pages 618–623. IEEE, 2017.
  13. Connah Kendrick, Bill Cassidy, Joseph M. Pappachan, Claire O’Shea, Cornelius J. Fernandez, Elias Chacko, Koshy Jacob, Neil D. Reeves, and Moi Hoon Yap. Translating clinical delineation of diabetic foot ulcers into machine interpretable segmentation, 2022.
  14. Manu Goyal, Neil D Reeves, Adrian K Davison, Satyan Rajbhandari, Jennifer Spragg, and Moi Hoon Yap. Dfunet: Convolutional neural networks for diabetic

- foot ulcer classification. *IEEE Transactions on Emerging Topics in Computational Intelligence*, 4(5):728–739, 2018.
15. Nora Al-Garaawi, Raja Ebsim, Abbas F.H. Alharan, and Moi Hoon Yap. Diabetic foot ulcer classification using mapped binary patterns and convolutional neural networks. *Computers in Biology and Medicine*, 140:105055, 2022.
  16. *dupeGuru*, 2018. [Online] Available from: <https://dupeguru.voltaicideas.net/> [Accessed: 7th June, 2022].