

Training a high-performance retinal foundation model with half-the-data and 400 times less compute

Justin Engelmann*, Miguel O. Bernabeu

*justin.engelmann@ed.ac.uk

JE: Centre for Medical Informatics, Usher Institute, University of Edinburgh, Edinburgh, UK; School of Informatics, University of Edinburgh, Edinburgh, UK

MB: Centre for Medical Informatics, Usher Institute, University of Edinburgh, Edinburgh, UK

Abstract

Artificial Intelligence holds tremendous potential in medicine, but is traditionally limited by the lack of massive datasets to train models on. Foundation models, pre-trained models that can be adapted to downstream tasks with small datasets, could alleviate this problem. Researchers at Moorfields Eye Hospital (MEH) proposed RETFound-MEH, a foundation model for retinal imaging that was trained on 900,000 images, including private hospital data. Recently, “data-efficient” DERETFound was proposed that provides comparable performance while being trained on only 150,000 images that are all publicly available. However, both these models required very substantial resources to train initially and are resource-intensive in downstream use. We propose a novel Token Reconstruction objective that we use to train RETFound-Green, a retinal foundation model trained using only 75,000 publicly available images and 400 times less compute. We estimate the cost of training RETFound-MEH and DERETFound at \$10,000 and \$14,000, respectively, while RETFound-Green could be trained for less than \$100, with equally reduced environmental impact. RETFound-Green is also far more efficient in downstream use: it can be downloaded 14 times faster, computes vector embeddings 2.7 times faster which then require 2.6 times less storage space. Despite this, RETFound-Green does not perform systematically worse. In fact, it performs best on 14 tasks, compared to six for DERETFound and two for RETFound-MEH. Our results suggest that RETFound-Green is a very efficient, high-performance retinal foundation model. We anticipate that our Token Reconstruction objective could be scaled up for even higher performance and be applied to other domains beyond retinal imaging.

Introduction

Artificial Intelligence has many promising applications in medicine, but the lack of large labelled datasets and access to vast computational resources present substantial bottlenecks [1]. Domain-specific “foundation models” that can be efficiently adapted to various downstream tasks are being proposed to remedy this issue [2], [3], [4]. Zhou et al. [5] from Moorfields Eye Hospital (MEH) recently proposed such a foundation model for retinal imaging, called “RETFound”. This could help unlock the potential of Artificial Intelligence in ophthalmology, where low-cost, non-invasive retinal colour fundus images are routinely used to screen for and

diagnose retinal disease, a key public health burden [6], [7]. Artificial Intelligence could aid in the interpretation of these images [8], [9], [10], [11], [12], [13]. These images capture the retina and its blood vessels in detail and thus allow inferences about the systemic health of individuals [14], [15], [16], [17], too, a field of study known as “oculomics” [18].

RETFound, henceforth referred to as “RETFound-MEH” for Moorfields Eye Hospital, presents a substantial contribution to the field, but took substantial resources to train: 900,000 retinal colour fundus images that are largely not publicly available and two weeks of eight high-end A100 datacentre-grade GPUs, which we estimate to have cost over \$10,000 (see Methods for calculations). Note that this is only for the training the final model and not additional experimentation that is typically necessary when training deep learning models. This level of resource consumption makes further scaling up expensive, both financially and environmentally [19], [20], and puts foundation model development out of reach of all but the most well-resourced labs.

More recently, Yen et al. proposed a data-efficient “DERETFound” [21] that was trained using only 150,000 images that are all publicly available. To accomplish this, they trained a diffusion model to generate over a million synthetic colour fundus images and then trained their model first on the synthetic and then on the real images. However, while it lowers the bar for dataset size, it does so at substantially increased computational cost due to training the generator and then using it to generate images. Overall, DERETFound required about 45% more computational resources than RETFound-MEH.

Both RETFound-MEH and DERETFound use the “Masked AutoEncoder” (MAE) [22] self-supervised learning strategy including the original hyperparameters, which Zhou et al. found to be more effective than other self-supervised strategies like SimCLR [23], SwAV [24], DINO [25] or MoCo-v3 [26]. However, these strategies are proposed in the general computer vision literature and thus designed for a task that differs substantially from the development of retinal image foundation model. First, they train models from scratch using randomly initialised weights, which is particularly challenging and computationally expensive for modern vision transformer architectures [27] that do not have strong inductive biases like traditional convolutional neural networks. Both RETFound-MEH and DERETFound use these pre-trained weights and thus only need to adapt the model to retinal imaging. Second, general computer vision uses natural images which have a different structure and more high-level diversity. In a dataset like ImageNet, an image of a dog is very different from an image of a car, which both in turn are very different from an image of a plate of food. In ophthalmology, a healthy eye and an eye with age-related macular degeneration might only differ through the presence of small deposits called drusen, which show up as tiny specks on the images. Third, datasets in general computer vision are very large with tens of millions [28] or even billions of images [29].

Design choices common in computer vision that were adopted by both RETFound-MEH and DERETFound include the MAE self-supervised learning approach, as well as a low resolution of 224 by 224 pixels, a fraction of the 3-4,000 pixels that modern retinal cameras tend to acquire,

and using the “large” variant of the vision transformer architecture [27], which makes them somewhat computationally expensive during inference, i.e. when processing images with the model for downstream tasks. While training costs are large, if models are adopted in routine care, the inference costs are recurring and thus a substantial factor with a non-negligible environmental footprint [19]. MAE involves pixel-level re-construction of images which is effective for the high-level diversity in natural images, not well suited for capturing small structures.

In this work, we propose a novel Token Reconstruction self-supervised learning strategy that focuses on higher-level, abstract features and is designed for making domain-specific foundation models instead of training models from scratch for general computer vision. See the Methods section and Figure 5 for a detailed explanation. We use our Token Reconstruction objective to train RETFound-Green, a high-performance retinal foundation model. Figure 1 gives an overview of how RETFound-Green compares with RETFound-MEH and DERETFound. Our strategy allows RETFound-Green to be trained far more efficiently in terms of data and compute, and yields a model that is substantially more lightweight in downstream applications while not being systematically worse than the previous models. On the contrary, in our experiments, RETFound-Green obtains 14 statistically significant wins, compared to six for DERETFound and 2 for RETFound-MEH. All of our experiments use open data and the methodology underlying each of the comparisons is explained in detail in the Methods section. We make RETFound-Green openly available and expect that it will not only be a useful tool for researchers in ophthalmic AI but democratise both access to and development of foundation models. Our Token Reconstruction objective is not explicitly designed for retinal image analysis and thus might find application in other medical and non-medical domains.

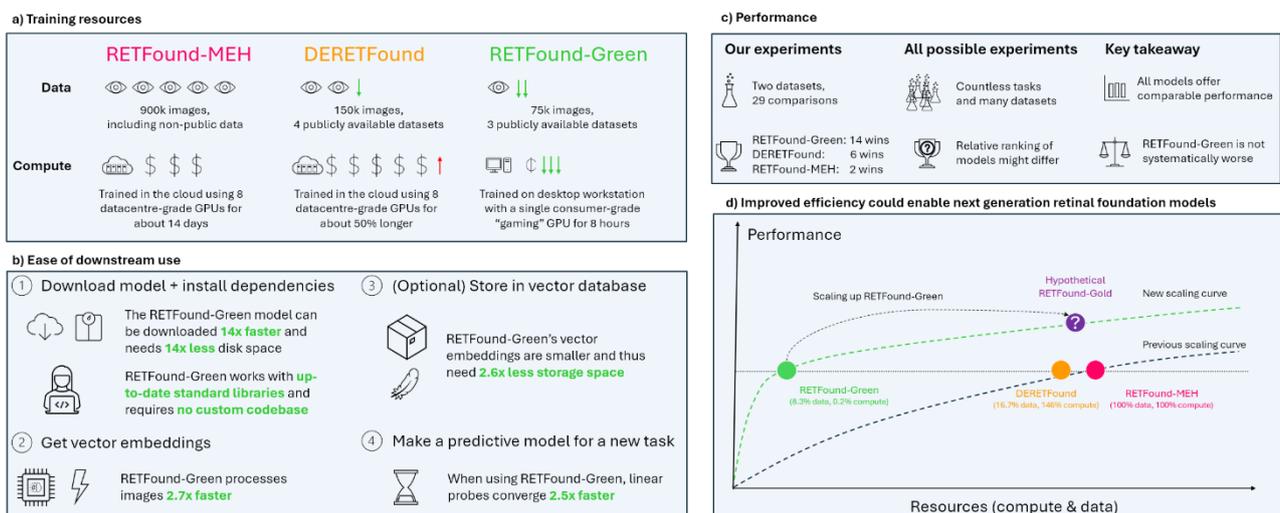


Figure 1: Comparison of RETFound-MEH, DERETFound, and RETFound-Green. a) RETFound-Green was trained with substantially less data and compute. b) RETFound-Green is easier and more efficient in downstream applications. c) RETFound-Green does not perform generally worse and better in most cases in our experiments, but this might differ across datasets and tasks. d) RETFound-Green is a far more efficient approach and could be scaled up to yield an even higher performance, next-gen foundation model.

Results

Training and downstream use efficiency

	RETFound-MEH	DERETFound	RETFound-Green	Green vs best
Training data	904,170 images	150,786 images	75,000 images	2x less
Training compute	112 A100 days	163 A100 days	~0.27 A100 days	>400x less
Training cost (monetary/carbon, estimate)	~\$10,000 / 81kg of coal burned	~\$14,000 / 117kg of coal burned	<\$100 / 0.2kg of coal burned	>100x less
Training hardware	8x top datacentre GPUs (total VRAM: 320 GB)	8x top datacentre GPUs (total VRAM: 640 GB)	1x top consumer gaming GPU (total VRAM: 24 GB)	>8x less
Disk space (model)	1.12 GB (our optimisation, 3.68 GB originally)	1.12 GB (our optimisation, 3.68 GB originally)	0.09 GB	14x less
Disk space (1 million embeddings)	39.1 GB	39.1 GB	14.6 GB	2.6x less
Inference speed (same hardware)	6 img / s	6 img / s	16 img / s	2.7x faster
Linear probe speed (same hardware)	2.45 s / task	2.40 s / task	0.96 s / task	2.5x faster
Performance (Only counting wins with p<0.05)	At least comparable (Various BRSET tasks [Fig. 2]: 5 wins for RETFound-Green, 4 wins for DERETFound, none for RETFound-MEH. Diabetic retinopathy grading [Fig. 3]: 9 wins for RETFound-Green, 2 wins each for DERETFound and RETFound-MEH.)			Not generally inferior

Table 1: Training resources and downstream use efficiency. See the following sections for detailed results regarding the performance on downstream tasks. See the Methods section for detailed explanations for all of the other rows.

RETFound-Green uses a novel Token Reconstruction self-supervised learning pre-training objective that allows it to be trained far more efficiently than RETFound-MEH and DERETFound which use the Masked Autoencoder (MAE) [22] objective. This allows for substantial improvements in efficiency as shown in Table 1.

RETFound-Green was trained with only 75,000 images, half of what DERETFound was trained on and 12 times less than the original RETFound-MEH. Furthermore, RETFound-Green required two orders of magnitude less computational resources, 400 times less than the original RETFound-MEH and 600 times less than DERETFound which required substantial compute resources for generating synthetic images. Thus, RETFound-Green was trained at a substantially lower cost. This also translates into a smaller estimated carbon footprint [20]: Training RETFound-MEH had an environmental impact comparable to burning 81kg of coal, DERETFound 117kg of coal, and RETFound-Green only 0.2kg of coal.

In addition to being more efficient to train, RETFound-Green is also more efficient in downstream use. The model is 14x smaller and can thus be stored and downloaded more easily, which especially benefits researchers with slow internet connections. Note, that this already factors in an improvement we made to RETFound-MEH and DERETFound that more than halves their file size without loss of performance, as described in detail in the Methods section. It also provides denser embeddings that require 2.6 times less space, which makes maintaining and sharing a vector database of image embeddings more efficient. Obtaining embeddings of images is about 2.6 times faster and thus more accessible even in lower resource settings. Finally, as the embeddings are denser, fitting a predictive model using those embeddings is also faster.

Another advantage that is harder to quantify but very important in practice is the ease-of-use and maintainability of the software code. In software development, there is a concept of “technical debt” where imperfect solutions lead to problems and increased workload in the future, similar to paying back a loan with interest. Being easier to set up and more maintainable is particularly important as users of medical foundation models might be less familiar with programming than those developing them. Thus, technical know-how can be a real bottleneck that is exacerbated by reliance on complex or outdated code.

Specifically, RETFound-MEH requires a five year old version of the Python programming language (version 3.7.5) which is now considered “end-of-life” and no longer officially supported [30]. Anecdotally, this meant that we were unable to use RETFound-MEH in a Scottish safe haven as this version of Python is considered a security risk by the IT team and thus could not be installed. RETFound-Green works with current, up-to-date versions of Python. Both RETFound-MEH and DERETFound further require a four year old version of the popular PyTorch Image Models library [31] and a modified version of the codebase for MAE [22] released by Meta two years ago which contains hundreds of lines of code. By contrast, RETFound-Green works with the most recent version of the PyTorch image models library and can be loaded with two lines of code.

Performance on diverse downstream tasks

We compare the performance of the three foundation models on nine different downstream tasks using BRSET [32], a Brazilian dataset of 16,266 colour fundus images. We chose this dataset as it is relatively large and richly annotated, which enables us to compare the models on a variety of downstream tasks, namely classifying abnormality of three key anatomical landmarks visible in colour fundus images (macula, optic disc, and retinal vessels), three different retinal diseases (age-related macular degeneration, macular edema, and diabetic retinopathy), and three non-retinal disease tasks (image quality, diabetes status, and insulin usage for diabetics).

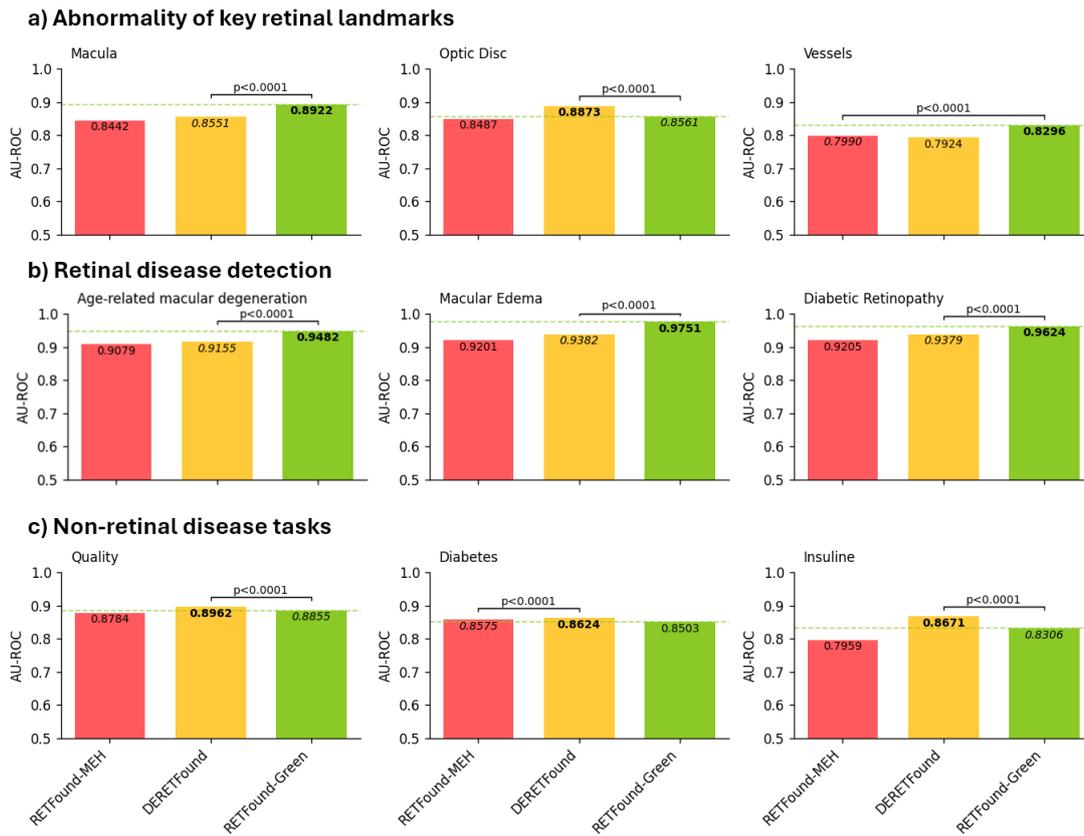


Figure 2: Performance for a variety of tasks in BRSET. The horizontal green dashed line indicates the performance of RETFound-Green to aid visual comparison. For robustness, reported results are the median of 100 bootstrap samples of the test set. Best result for each task in bold, the bar with p-value indicates the result of a Wilcoxon signed-rank test between the best and second best methods across the 100 bootstrap samples.

The results are shown in Figure 2. Counting only statistically significant wins with $p < 0.05$, RETFound-Green performs best for five of the nine tasks, while DERETFound wins four times. For most tasks, the best method wins by a reasonably large margin, except for Quality and Diabetes where the differences are statistically significant but all three methods obtain very comparable performance. Overall, RETFound-Green obtains the most wins and generally performance is at least comparable across all tasks but for Optic Disc abnormality and Insulin usage, DERETFound obtains substantially better results.

Performance on fine-grained diabetic retinopathy grading

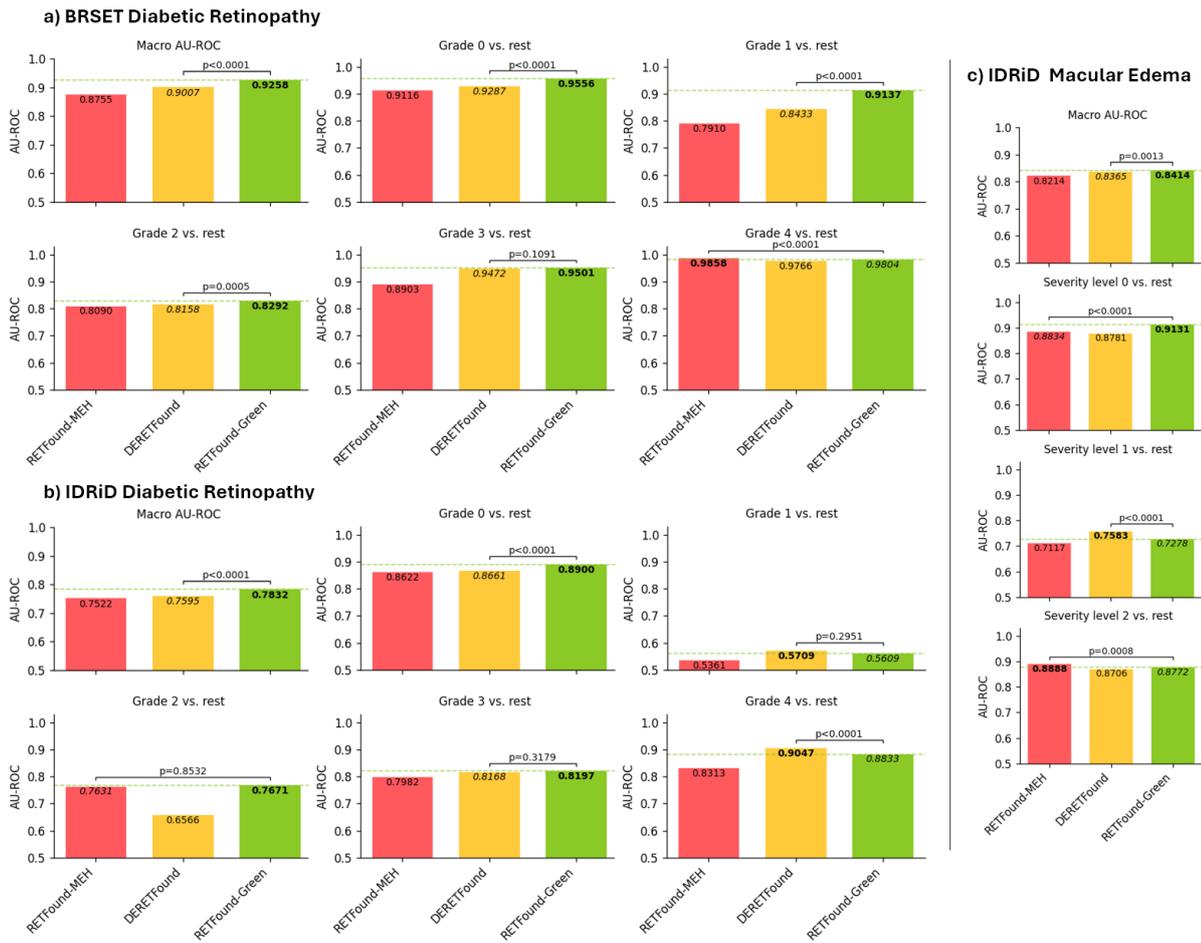


Figure 3: Performance for fine-grained diabetic retinopathy grading. The horizontal green dashed line indicates the performance of RETFound-Green to aid visual comparison. For robustness, reported results are the median of 100 bootstrap samples of the test set. Best result for each task in bold, the bar with p-value indicates the result of a Wilcoxon signed-rank test between the best and second best methods across the 100 bootstrap samples.

In addition to the variety of different tasks in the previous section, we further look at fine-grained diabetic retinopathy related tasks. Diabetic retinopathy is a very common retinal disease and screening for it as well as grading its severity is key part of ophthalmic care, thus it makes a good test case for detailed retinal disease assessment. Here, we use the severity grades provided in BRSET and further consider the Indian Diabetic Retinopathy Image Dataset (IDRiD) [33] that provides severity grades for diabetic retinopathy and macular edema, which is common in patients with diabetic retinopathy.

Figure 3 shows the results. RETFound-Green obtains nine statistically significant wins, whereas RETFound-MEH and DERETFound each win in two cases. In three cases, namely diabetic retinopathy grades 1, 2, and 3 for IDRiD, the best method did not outperform the second best one in a statistically significant way. However, it is worth noting that in all three cases, RETFound-Green was either the winner or runner-up, while the third method performed substantially worse than the first two. Thus, overall RETFound-Green appears to offer the best performance for fine-grained disease grading tasks.

Effectiveness of the RETFound-Green Token Reconstruction pre-training objective

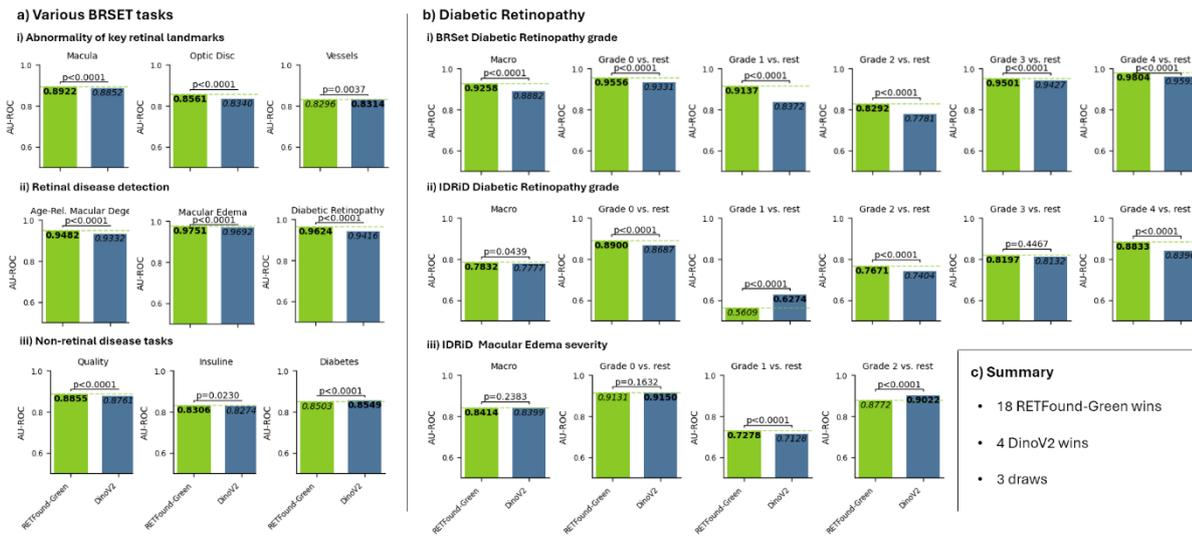


Figure 4: Comparison between RETFound-Green and DinoV2. The horizontal green dashed line indicates the performance of RETFound-Green to aid visual comparison. For robustness, reported results are the median of 100 bootstrap samples of the test set. Best result for each task in bold, the bar with p-value indicates the result of a Wilcoxon signed-rank test between the best and second best methods across the 100 bootstrap samples.

To evaluate the effectiveness of our Token Reconstruction pre-training object, we compare RETFound-Green with the original DinoV2 [34] model that we initialised our weights with across all tasks from the previous two sections. Figure 4 shows the results of this comparison. RETFound-Green obtains 18 statistically significant wins, while DinoV2 obtains 4 wins, and there are 3 draws where neither method performed significantly better. These results suggest that while DinoV2 already presents a very strong baseline, our Token Reconstruction pre-training objective was indeed effective in pre-training RETFound-Green for retinal images.

Discussion

This work presents RETFound-Green, a new high-performance retinal image foundation model that was trained using a novel Token Reconstruction pre-training objective. This novel objective allowed us to train RETFound-Green with 50% less data and 400x less computational resources compared to the “best of both worlds” of the two previously proposed models, namely DERETFound’s data efficiency and RETFound-MEHs compute efficiency. RETFound-Green is also far more efficient and easier to use in downstream applications. Despite these substantial improvements in efficiency, RETFound-Green does not perform systematically worse and in fact in our experiments it has 14 statistically significant wins compared to four for DERETFound and two for RETFound-MEH.

While AI-powered advancements in healthcare hold great promise, there is a risk that they could lead to new or exacerbated disparities as the resources necessary to develop or make use of them are not accessible to everyone [35]. We believe that RETFound-Green will not only serve as a foundation model for retinal image analysis but also democratize foundation models in multiple ways: First, democratising access to foundation models. RETFound-Green is far

more efficient in downstream applications, which will especially benefit researchers with comparatively fewer resources. For example, on a 128Kb/s internet connection, it would take 2 days and 15 hours to download RETFound-MEH or DERETFound as shared by the authors, 19.5 hours with our optimisation of their models, but RETFound-Green could be downloaded in just 1.5 hours. As RETFound-Green is 2.7 times faster when it comes to calculating embeddings, it could take a workload that previously would have taken a whole day down to less than nine hours, or allow such a workload to be run in the same time on hardware that is 2.7 times slower. If models are used routinely in healthcare, this can produce substantial emissions [19] and increased inference efficiency translates into more sustainable medical AI.

Second, RETFound-Green and the Token Reconstruction objective we introduce democratise foundation model development by requiring far less data than before: half of the data-efficient DERETFound and 12 times less than the original RETFound-MEH. This not only allows researchers without access to massive databases to develop their own foundation models, but could also enable future foundation models to be substantially fairer than current ones. As the required dataset size is now less than 100,000 images, it would be feasible to collect a comparatively small yet representative and unbiased dataset that includes data from diverse patients in terms of ethnicity, sex, disease status, etc. and from a variety of healthcare contexts across the world.

Third, the reduced computational cost of our pre-training objective democratises foundational models by making it accessible to researchers without substantial financial and computational resources, as well as by making it more sustainable by reducing the environmental impact which makes it fairer towards future generations and current generations in developing countries, who are most affected by the impact of climate change. Concretely, the previous approaches required over a hundred days of high-end datacentre GPU days that would cost thousands of dollars, whereas our model was trained for 8 hours on a consumer-grade GPU that could be found in a high-end “gaming computer”. By increasing the training time to a day or two, it would be possible to adapt our approach to instead also run on a modern low-end gaming computer, which would be within reach of many researchers globally. The massively reduced environmental impact means that future and scaled-up foundation models can be developed in good conscience.

It is important to note that our approach, including the Token Reconstruction objective, are in no way specific to retinal imaging and likely would work similarly well for developing foundation models for other imaging modalities, both medical and non-medical. We focus on retinal imaging is simply because we are familiar with this domain. Given the order-of-magnitude gains in efficiency we observe, we think it is imperative that we share our results promptly to hopefully avoid new foundation models being trained with less efficient methods which would be highly wasteful based on our results. Furthermore, RETFound-Green already offers many practical advantages for researchers working on retinal image analysis. Thus, we chose not to delay

disseminating our results, instead of demonstrating effectiveness across a range of modalities first.

While our Token Reconstruction objective appears to be very effective, we do not claim that this is already the best possible strategy for developing foundation models. In fact, we did not tune or iterate on our approach at all, thus it would be surprising if there was no room left for improvement. However, at present our approach is already vastly more efficient than the approach used by RETFound-MEH and DERETFound. One important general lesson from this is that when developing models for specific domains, it is likely very suboptimal to simply copy the approach taken by researchers in general, natural image computer vision. The approaches are developed for datasets that are orders of magnitude larger and to train models entirely “from scratch”, i.e. using randomly initialised parameters.

RETFound-Green defies conventional wisdom that large amounts of data and compute are necessary for the current levels of performance by achieving the same level of performance with much less resources. This raises the possibility that a better version of RETFound-Green could be trained with a similar amount of resources to what was used for RETFound-MEH and DERETFound.

There are two additional potential benefits of RETFound-Green that we do not investigate in the current manuscript but that are important to note. First, the smaller model size would be a substantial benefit for federated learning, where the data remains distributed across multiple sites (e.g. different hospitals) and not centrally aggregated. A big bottleneck for federated learning is that either model parameters or gradient updates need to be sent back and forth at each iteration, many thousands or even hundreds of thousand times. For RETFound-MEH and DERETFound, this would require over 1GB of data transfer for each iteration, whereas for RETFound-Green it would only be 0.09 GB, which would result in massively reduced training time and costs. Second, RETFound-Green might be more privacy-preserving. The MAE pre-training objective involves reconstructing the exact pixels of the input image, whereas our Token Reconstruction objective focuses only on more abstract features. The RETFound-Green embeddings are also denser – 384 vs 1,024 numbers per image – which substantially reduces the risk of the embeddings allowing specific patients to be identified or images to be reconstructed.

While our results are overall strong, there are a number of limitations for this work. First, our evaluations focused on two datasets, one from Brazil and one from India. While we compare the models across a variety of tasks and find that RETFound-Green generally performs best, it is possible that for other datasets and tasks the relative performance of the three evaluated models could be different. At present, the most important insight is that RETFound-Green does not perform systematically worse while being far more efficient, which is well supported by the data. Still, in future work it would be interesting to compare the models across many datasets and an even broader selection of tasks. Second, although achieving the same level of performance with far fewer resources suggests that even better performance could be achieved

when scaling our approach up, we do not attempt this in the current manuscript as we do not have access to that level of resources ourselves. This is a very interesting possibility that should be explored by future work. Third, while our approach is not specific to retinal images and likely could be applied to many different imaging domains, we likewise do not investigate this in the current manuscript. Finally, our approach being generic rather than specific to retinal images is also a weakness as adding elements specific to retinal image analysis and ophthalmology, such as integrating additional information about the patient and their symptoms or considering longitudinal images, could further improve the utility of our model for practical applications, which we plan to investigate in future work.

In conclusion, we present RETFound-Green trained with a novel Token Reconstruction approach that allows us to match the performance of previous retinal image foundation models with far less resources while being substantially more efficient and easy to use in practice. This is useful for researchers in the field and could democratise both access to and development of foundation models in retinal imaging and beyond.

Methods

RETFound-Green

Like RETFound-MEH and DERETFound, RETFound-Green uses the original vision transformer architecture [27] but with four extra “register tokens” [36]. This is a minor and straightforward modification of the architecture that has been observed to yield smoother attention maps. More recent and advanced vision transformer architectures could be explored in the future, for example using group tokens [37] or query pooling [38] to reduce the internal dimensionality. In this work, we chose a similar architecture to make the comparison more direct. RETFound-MEH and DERETFound use the “large” version of the vision transformer, RETFound-Green “small” version. There is also a “base” version, so the small one is two steps down and much smaller. This allows RETFound-Green to process the images at higher resolution of 392x392 instead of 224x224 used by the other two models while still being more efficient at the same time.

RETFound-MEH and DERETFound start with pre-trained weights from the Masked Autoencoder (MAE) paper [22] and then train these weights on retinal images using the MAE strategy and codebase. We start with weights from DinoV2 [34], a self-supervised learning method like MAE. However, unlike the other two foundation models, we only use the DinoV2 weights and then train those on retinal images using our own Token Reconstruction strategy, described in the next section, and our own codebase.

Self-supervised training via Token Reconstruction

We propose a novel method for self-supervised training of pre-trained models. We note that MAE [22] and DinoV2 [34] were devised for training models entirely from scratch, i.e. using randomly initialised weights. These strategies are very computationally expensive and require vast amounts of data. When developing foundation models for medical imaging, we can start

with models that were already pre-trained on natural images. In fact, this is the same approach as RETFound-MEH and DERETFound use. However, they not only use the pre-trained MAE weights as a starting point but the MAE self-supervised learning objective, including the hyperparameters. This might not be optimal as the MAE strategy and the hyperparameters were chosen for training from scratch on natural images, which is complex and resource-intensive, especially for vision transformers.

When adapting the existing model to retinal imaging, we have two objectives for our self-supervised training approach. First, instilling knowledge about the general appearance and structure of retinal images. Second, optimising the features of our model for this space without unlearning existing, useful features. We propose a novel token reconstruction strategy that is straight-forward yet effective: We take a frozen version of the model we are adapting, pass retinal images through it and obtain the output tokens. Our model, a second non-frozen copy, is given noisy versions of the same images as input and its output tokens are compared to the first model’s output tokens, and our model incurs a loss for how different the two are. In other words, our model has to reconstruct the output tokens from noisy inputs.

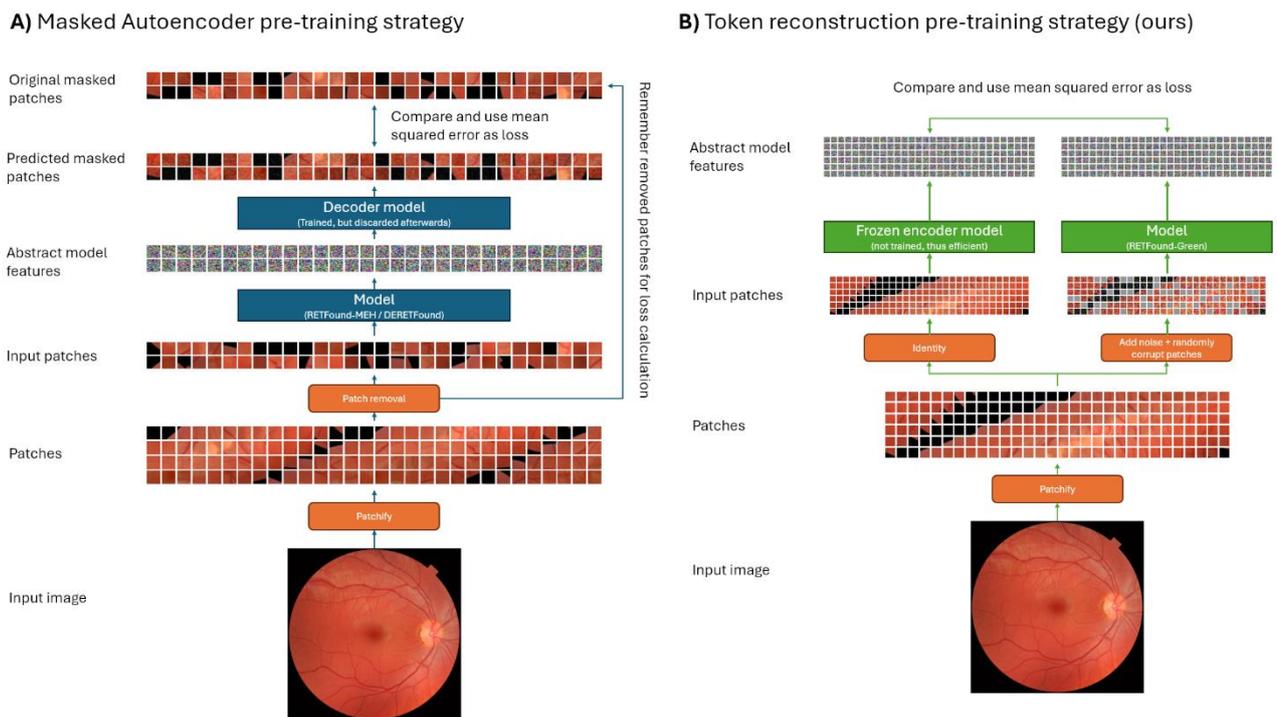


Figure 5: Comparison between the Masked Autoencoder (MAE) strategy (a) used by RETFound-MEH and DERETFound with our proposed Token Reconstruction strategy (b) that we used to train RETFound-Green. Our strategy focuses on the abstract features that are used in downstream applications, whereas MAE focuses on exact pixel values. MAE requires training a decoder model which increases computational cost but this model is discarded after training. Our strategy also uses an auxiliary model, but this is kept frozen and thus gradients do not need to be computed which lowers computational cost substantially.

The intuition behind the Token Reconstruction objective is that to match the output of the frozen model, our model needs to use information from the non-corrupted input tokens to predict the correct outputs for all tokens. This requires understanding the structure and co-dependencies

of retinal images. For instance, if a corrupted patch contains lesions, then it is likely that the rest of the image either also contains lesions or other anatomical features that make the presence of lesions more probable. Furthermore, some anatomical structures occur in specific locations such as the optic disc or macula. To reconstruct these, the model benefits from learning what structures tend to occur in which locations. Likewise, the pixel-wise noise needs to be ignored by our model, which requires understanding what is part of the original image and what is added noise.

We use two simple noise strategies: First, a random subset of input tokens is replaced with a trainable “corruption token” which effectively erases the information about a patch of the input image. Second, we apply pixel-wise Gaussian noise to the input image. The detailed parameter settings are described in the section on training parameters below. Our strategy is straightforward to implement and computationally efficient. We note that there are many possible extensions of approach that could be tried, for example additional strategies of making the image noisy, or reconstructing the outputs of multiple different frozen models at the same time.

Datasets for pre-training and downstream evaluations

We used three publicly available datasets for pre-training: AIROGS [39], DDR [40], and ODIR-2019. We used all images in DDR and ODIR-2019 and then selected a random sample of 53,327 from AIROGS to achieve our desired target amount of exactly 75,000 images. We make the subset of images used available on our GitHub to aid reproducibility.

For our downstream evaluations, we used the Brazilian Multilabel Ophthalmological Dataset (“BRSET”) [32] which contains very rich annotations that enable us to compare models across a variety of tasks. We randomly split BRSET into training (80%) and testing (20%) sets at the patient level, ensuring that no patient is present in both sets to avoid data leakage. These splits are also made available on our GitHub. We further used the Indian Diabetic Retinopathy Image Dataset (“IDRiD”) [33] where we use the official train-test-split.

Before applying our data pipeline during training, we resize all images to a resolution of 1024x1024 and save them with very high JPG quality. Colour fundus images are circular and thus the images should be square. However, in some datasets there is considerable horizontal black space. Naïve resizing to a square resolution, as is commonly done in the literature, does not preserve the aspect ratio and leads to the images being squashed. Some images are cropped slightly at the top and bottom, which also gives them a non-square aspect ratio. We detect the fundus area and then crop horizontal black space, or pad the top and bottom, to ensure the images are properly cropped and square before resizing. We implement this because we think it is a principled way to resizing retinal images, but we do not think that this has any substantial impact at all on performance.

For downstream evaluations, we simply resize the images to the desired size and normalise by each model’s specific normalisation constants. We intentionally do not use our more complex

preprocessing on the downstream tasks to make the comparisons fair by ensuring that better preprocessing plays no role.

RETFound-Green training parameters

RETFound-Green was trained for 120 epochs, batch size of 128, using AdamW [41] with a maximum learning rate of $5 * 10^{-5}$, betas $\beta_1 = 0.9, \beta_2 = 0.99$ and weight decay $5 * 10^{-4}$. No weight decay is applied to the bias parameters, a common practice in modern deep learning. We use a cosine learning rate scheduler [42] with a warmup of 10 epochs where the learning rate linearly increases from its minimum $5 * 10^{-9}$ to its maximum and a cool-down of 20 epochs where it is kept at its minimum. We use automatic mixed precision with bfloat16, a “half-precision” data type that uses 16 instead of standard 32-bit floating point numbers that is optimised for deep learning. This reduces memory consumption and speeds up training. The maximum gradient norm is clipped to 0.1.

For the Token Reconstruction objective, we use a corruption ratio sampled from $U(0, \frac{1}{3})$, a pixel corruption with noise sampled from $N(0, 0.2)$. The last image in each batch is kept uncorrupted so the model is robust to the absence of corruption tokens. As loss function we use simple mean squared error. For the projection, we use a small residual MLP consisting of LayerNorm, a linear layer, a GELU activation, and another linear layer. Each of the linear layers had a dimension of 384, same as the dimensionality of our model. The projector takes the model’s final representation as input, and the projector’s output is then scaled by an element-wise learnable parameter and added to the final representations before loss computation.

We apply slight colour jitter and rotations 25% of the time, pad the sides of the images by random value between 33 and 150 pixels to simulate poorly cropped images 10% of the time, and scale the image with a random value between 80 and 120%. We then use one of three ways to obtain an image of our target resolution of 392x392: standard resizing, or a random resized crop with a scale of 70-100%, or a centre crop after scaling the image up by 30%. Note that first properly cropping the images and then simulating poor cropping randomly, we decorrelate poor cropping from specific datasets, which might make our model more robust.

RETFound-Green uses simple normalisation with mean and standard deviation parameters of 0.5 for all three colour channels, whereas RETFound-MEH and DERETFound use the statistics of the ImageNet dataset of means 0.485, 0.456, 0.406 for red, green, and blue channels, and standard deviations of 0.229, 0.224, 0.225. The model learns to adjust to these constants during training, so this is a very small optimisation to increase convenience in downstream use as users of RETFound-Green only need to remember 0.5, instead of looking up these values.

We use identical settings for the version at a lower resolution of 224x224, noting that efficiency likely could be improved by increasing the batch size and scaling down the learning rate accordingly. However, we did not tune these parameters for RETFound-Green nor the lower resolution version. Instead, we chose reasonable values that worked well out of the box.

Predictive modelling for downstream tasks

We obtain vector embeddings for the images from each of the foundation model and then fit a linear model with logistic linkage function to the downstream training set, also known as “linear probing” in the machine learning community. We use default values in the popular scikit-learn [43] machine learning library, except for increasing the maximum fitting iterations to 20,000 to ensure that all models converge successfully. Note that in machine learning, unlike traditional statistics, a L2-weight penalty is used by default which substantially improves convergence and tends to increase predictive performance. We standardise the data to 0 mean and unit variance using scikit-learn’s StandardScaler to help convergence and ensure the weight penalty has the same effect for all variables. We follow machine learning best practices and estimate the parameters for standardisation only on the training set, which avoids data leakage and simulates the scenario where the test data is not available at training time.

This strategy for adapting foundation models has many advantages over fine-tuning the whole model. First, it can be done using only vector embeddings and labels without requiring the images themselves. This allows developing new models just from vector databases. Second, it is very computationally efficient and can be done on low-end hardware, especially if the vector embeddings are already computed. But even if not, this strategy only requires inferences on each image once (also known as a “forward pass” in machine learning), whereas fine-tuning typically requires multiple passes over each image (once per epoch) and each of those passes is computationally more expensive as in addition to the forward pass we also need to compute gradients and take optimisation steps (also known as “backward pass”). Third, it requires much less technical know-how than fine-tuning a deep learning model and could be done in statistical programming languages like R that clinicians are more likely to be familiar with as well as Python, whereas deep learning is best done in languages like Python or C++.

Computational resources

We measure compute in terms of “A100 days”, which means the equivalent of training on a single Nvidia A100 Tensor Core datacentre GPU for 1 day. The original RETFound was trained on 8 A100s for two weeks, so $8 \times 14 = 112$ A100 days. DERETFound used 2 A100 days for training the diffusion model, about 113 days for generating images with that model, and finally 8 A100s for 6 days for training the DERETFound model itself, so $2 + 113 + (6 \times 8) = 163$ A100 days. RETFound-Green was trained for 8 hours on a desktop computer with a single Nvidia RTX 4090 GPU, a consumer-grade “gaming” GPU card. A 4090 is roughly equivalent to 0.82 the performance of the slowest A100, the 40GB cloud version, in the lambda labs benchmarks [44], or 0.63 compared the faster A100 80GB. We take the multiplier that is least favourable to RETFound-Green and estimate that it used $(8/24) \times 0.82 = 0.27$ A100 days.

For the estimated training costs, we use the on-demand price for a 8x A100 40GB cloud machine on Amazon Web Services of \$32.77 per hour (<https://aws.amazon.com/ec2/instance-types/p4/>) and round it down to \$30. Prices differ across providers and regions, and change over time, so these are ballpark estimates. DERETFound uses the 80GB variants for pre-training

which are more expensive but we use the same, lower price of the 40GB variant for all methods. We take the price per day and multiply it by the estimated A100 days divided by 8, as each machine has 8x A100. This gives \$10,080 for RETFound, \$14,670 for DERETFound, and \$24.30 for RETFound-Green. We round the other two methods down and RETFound-Green up to “<\$100”.

For estimating the CO₂ released, we use the Machine Learning Emissions Calculator (<https://mlco2.github.io/impact/>) proposed by previous work [20]. For a fair comparison, we use the same cloud provider and region, Amazon Web Services in US West (N. California), a rough mid-point between the UK and China, the two countries the compared foundation models were developed in. We think that this is fairer than taking into account the local energy mix, as we want to compare methods and not favour researchers based on where they happen to be located. For example, in Scotland where we are based, renewables provided 113% of the electricity consumption in January 2024 [45] and our model was trained overnight, when energy consumption tends to be low. Thus, we likely did not lead to any CO₂ being released, but this is just due to our location, not due to our methods. Thus, using the calculator and same cloud provider and region, we estimate RETFound-MEH to have generated 161kg CO₂ equivalent or 81kg of coal burned, DERETFound 234kg CO₂ equivalent or 117kg of coal burned, and RETFound-Green 0.39kg CO₂ equivalent or 0.2kg of coal burned.

The speed at which the vector embeddings are computed is estimated on the same hardware, a low-end workstation with a 12th Gen Intel i5-12600k CPU and an Nvidia RTX3060ti GPU, a last generation, low-end gaming card. This level of hardware is very accessible, even in comparatively low resource settings. For equal comparison, we measure inference speed using GPU-acceleration with full float32 precision and a batch size of 1, following the example script released by Zhou et al.

(https://github.com/rmaphoh/RETFound_MAE/blob/main/RETFound_Feature.ipynb). Batching and mixed precision could further improve inference time. The same CPU is used for fitting of the linear probes.

Storage space for model weights and vector embeddings

Users of foundation models incur storage costs for the model itself, and additionally for each vector embedding. For the model weights for RETFound and DERETFound take up 3.68GB each as shared by their original authors. However, we observe that this includes the image decoder and optimizer states, both of which are only used during training and not necessary in downstream use. The optimizer states are especially expensive as they contain one value for each parameter of the foundation model. By removing both of those and only leaving the weights of the foundation models themselves, we can optimise their file size to 1.12GB, a reduction of 58%. We use this optimised version for comparison to be maximally fair. RETFound-Green takes up 83.8MB, which we round up to 0.09GB.

RETFound and DERETFound have a vector embedding feature dimension of 1,024, RETFound-Green's dimension is only 384 and thus its vector embeddings require 2.67x less storage space. Storing 1 million embeddings would take about 39.1GB for RETFound and DERETFound, and 14.6GB for RETFound-Green without any compression, storing standard 32-bit floats.

Evaluation metrics and statistical analysis

We focus on the Area Under the Receiver Operator Characteristic curve (AU-ROC), a widely used ranking metric that summarises sensitivity and specificity across all possible decision thresholds. That makes the AU-ROC preferable for general comparisons over metrics that require binarization of predictions as it captures more information.

A common claim in the literature is that for datasets with class imbalance, the Area Under the Precision Recall curve (AU-PR) should be preferred. However, this is not well supported by evidence, as discussed in recent work by McDermott et al. [46] who further show theoretically and empirically that AU-ROC is robust to class imbalance, while selecting models using AU-PR can lead to disparities across subpopulations. In their conclusion they explicitly state that the AU-ROC might be desirable in domains like healthcare, while AU-PR might not be reliable "in settings where equity and fairness are imperative". Thus, we focus on the AU-ROC in this work.

To ensure robustness of metrics and calculate p-values for model comparisons, we first bootstrap the test set 100 times, compute the AU-ROC for each sample, and then report the median value. This captures uncertainty over different test set distributions and ensure that our metrics are not unduly influenced a few individual samples. We note that Logistic Regression with weight penalty has a unique solution up to tiny differences due to floating point precision and is thus deterministic, so there is no uncertainty in the model fitting procedure.

We then do a Wilcoxon signed-rank test across all 100 bootstrap AUC values between the best and second-best methods to test whether the two methods have non-equal performance [47], [48]. We use the same 100 test set bootstrap samples for both methods so we can do a paired comparison. The Wilcoxon test is a non-parametric alternative to a paired t-test.

Statistical comparisons of classifiers are a complex topic and frequently subject to mistakes, for example Wilcoxon is appropriate for classifier comparisons whereas a t-test is not [48]. Likewise, we intentionally do not provide confidence intervals here to avoid misinterpretation. First, for non-normal data they tend to provide incorrect coverage [49]. Second, and more crucially, the performance across the bootstrap samples for different methods are non-independent as the variance in performance is driven by the sampling. Thus, overlapping confidence intervals do not imply non-significant differences. If method A performs better than method B for all or a great many of the individual bootstrap samples, then we should conclude that A is better than B, even if A's advantage over B is small relative to the variance between bootstrap samples. This is what the Wilcoxon test captures.

Model and code availability

We share the trained RETFound-Green model, as well as model, training and evaluation code on our GitHub: https://github.com/justinengelmann/RETFound_Green

Data availability

All datasets used in this manuscript are publicly available. We further make our code, model, and additional information to aid reproducibility available on our GitHub.

AIROGS dataset: <https://zenodo.org/records/5793241>

DDR dataset: <https://github.com/nkicsl/DDR-dataset>

ODIR-2019 dataset (registration required): <https://odir2019.grand-challenge.org/Download/>

BRSET dataset (registration required): <https://physionet.org/content/brazilian-ophthalmological/1.0.0/>

IDRiD dataset: <https://ieee-dataport.org/open-access/indian-diabetic-retinopathy-image-dataset-idrid>

Acknowledgements

We thank the authors of the original RETFound, Dr Yukun Zhou, Prof. Pearse Keane and colleagues, as well as the authors of DERETFound, Prof. Yan Bo and colleagues, for their contribution to the field and particularly for making their models openly available which enables the comparisons in this work. We further thank the researchers that made the datasets used in this work available and the individuals who contributed their data to biomedical research.

JE was supported by the United Kingdom Research and Innovation (grant EP/S02431X/1), UKRI Centre for Doctoral Training in Biomedical AI at the University of Edinburgh, School of Informatics.

M.O.B. was supported by: Fondation Leducq Transatlantic Network of Excellence (17 CVD 03); EPSRC grant no. EP/X025705/1; British Heart Foundation and The Alan Turing Institute Cardiovascular Data Science Award (C-10180357); Diabetes UK (20/0006221); Fight for Sight (5137/5138); the SCONE projects funded by Chief Scientist Office, Edinburgh & Lothians Health Foundation, Sight Scotland, the Royal College of Surgeons of Edinburgh, the RS Macdonald Charitable Trust, and Fight For Sight.

References

- [1] P. Rajpurkar, E. Chen, O. Banerjee, and E. J. Topol, 'AI in health and medicine', *Nat Med*, vol. 28, no. 1, pp. 31–38, Jan. 2022, doi: 10.1038/s41591-021-01614-0.
- [2] A. B. Sellergren *et al.*, 'Simplified Transfer Learning for Chest Radiography Models Using Less Data', *Radiology*, vol. 305, no. 2, pp. 454–465, Nov. 2022, doi: 10.1148/radiol.212482.

- [3] R. J. Chen *et al.*, 'Towards a general-purpose foundation model for computational pathology', *Nat Med*, vol. 30, no. 3, pp. 850–862, Mar. 2024, doi: 10.1038/s41591-024-02857-3.
- [4] Z. Huang, F. Bianchi, M. Yuksekgonul, T. J. Montine, and J. Zou, 'A visual–language foundation model for pathology image analysis using medical Twitter', *Nat Med*, vol. 29, no. 9, pp. 2307–2316, Sep. 2023, doi: 10.1038/s41591-023-02504-3.
- [5] Y. Zhou *et al.*, 'A foundation model for generalizable disease detection from retinal images', *Nature*, vol. 622, no. 7981, pp. 156–163, Oct. 2023, doi: 10.1038/s41586-023-06555-x.
- [6] L. Pezzullo, J. Streatfeild, P. Simkiss, and D. Shickle, 'The economic impact of sight loss and blindness in the UK adult population', *BMC Health Services Research*, vol. 18, no. 1, p. 63, Jan. 2018, doi: 10.1186/s12913-018-2836-0.
- [7] G. C. Brown, 'Vision and quality-of-life.', *Trans Am Ophthalmol Soc*, vol. 97, pp. 473–511, 1999.
- [8] J. De Fauw *et al.*, 'Clinically applicable deep learning for diagnosis and referral in retinal disease', *Nat Med*, vol. 24, no. 9, pp. 1342–1350, Sep. 2018, doi: 10.1038/s41591-018-0107-6.
- [9] X. Liu *et al.*, 'A comparison of deep learning performance against health-care professionals in detecting diseases from medical imaging: a systematic review and meta-analysis', *The Lancet Digital Health*, vol. 1, no. 6, pp. e271–e297, Oct. 2019, doi: 10.1016/S2589-7500(19)30123-2.
- [10] D. S. W. Ting *et al.*, 'Artificial intelligence and deep learning in ophthalmology', *British Journal of Ophthalmology*, vol. 103, no. 2, pp. 167–175, Feb. 2019, doi: 10.1136/bjophthalmol-2018-313173.
- [11] L. Dai *et al.*, 'A deep learning system for predicting time to progression of diabetic retinopathy', *Nat Med*, vol. 30, no. 2, pp. 584–594, Feb. 2024, doi: 10.1038/s41591-023-02702-z.
- [12] J. Engelmann, A. D. McTrusty, I. J. C. MacCormick, E. Pead, A. Storkey, and M. O. Bernabeu, 'Detecting multiple retinal diseases in ultra-widefield fundus imaging and data-driven identification of informative regions with deep learning', *Nat Mach Intell*, vol. 4, no. 12, pp. 1143–1154, Dec. 2022, doi: 10.1038/s42256-022-00566-5.
- [13] A. D. Fleming *et al.*, 'Deep learning detection of diabetic retinopathy in Scotland's diabetic eye screening programme', *British Journal of Ophthalmology*, Sep. 2023, doi: 10.1136/bjo-2023-323395.
- [14] R. Poplin *et al.*, 'Prediction of cardiovascular risk factors from retinal fundus photographs via deep learning', *Nat Biomed Eng*, vol. 2, no. 3, pp. 158–164, Mar. 2018, doi: 10.1038/s41551-018-0195-0.
- [15] S. M. Zekavat *et al.*, 'Deep learning of the retina enables phenome- and genome-wide analyses of the microvasculature', *Circulation*, vol. 145, no. 2, pp. 134–150, 2022.
- [16] A. Villaplana-Velasco *et al.*, 'Fine-mapping of retinal vascular complexity loci identifies Notch regulation as a shared mechanism with myocardial infarction outcomes', *Commun Biol*, vol. 6, no. 1, pp. 1–13, May 2023, doi: 10.1038/s42003-023-04836-9.
- [17] R. Luben *et al.*, 'Retinal fractal dimension in prevalent dementia: The AlzEye Study', *Investigative Ophthalmology & Visual Science*, vol. 63, no. 7, pp. 4440–F0119–4440–F0119, 2022.
- [18] S. K. Wagner *et al.*, 'Insights into systemic disease through retinal imaging-based oculomics', *Translational vision science & technology*, vol. 9, no. 2, pp. 6–6, 2020.

- [19] A. V. Sadr *et al.*, ‘Operational greenhouse-gas emissions of deep learning in digital pathology: a modelling study’, *The Lancet Digital Health*, vol. 6, no. 1, pp. e58–e69, Jan. 2024, doi: 10.1016/S2589-7500(23)00219-4.
- [20] A. Lacoste, A. Luccioni, V. Schmidt, and T. Dandres, ‘Quantifying the Carbon Emissions of Machine Learning’. arXiv, Nov. 04, 2019. doi: 10.48550/arXiv.1910.09700.
- [21] B. Yan *et al.*, *Expertise-informed Generative AI Enables Ultra-High Data Efficiency for Building Generalist Medical Foundation Model*. 2024. doi: 10.21203/rs.3.rs-3766549/v1.
- [22] K. He, X. Chen, S. Xie, Y. Li, P. Dollár, and R. Girshick, ‘Masked Autoencoders Are Scalable Vision Learners’, arXiv.org. Accessed: Feb. 20, 2024. [Online]. Available: <https://arxiv.org/abs/2111.06377v3>
- [23] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, ‘A Simple Framework for Contrastive Learning of Visual Representations’. arXiv, Jun. 30, 2020. doi: 10.48550/arXiv.2002.05709.
- [24] M. Caron, I. Misra, J. Mairal, P. Goyal, P. Bojanowski, and A. Joulin, ‘Unsupervised Learning of Visual Features by Contrasting Cluster Assignments’. arXiv, Jan. 08, 2021. doi: 10.48550/arXiv.2006.09882.
- [25] M. Caron *et al.*, ‘Emerging Properties in Self-Supervised Vision Transformers’. arXiv, May 24, 2021. doi: 10.48550/arXiv.2104.14294.
- [26] X. Chen, S. Xie, and K. He, ‘An Empirical Study of Training Self-Supervised Vision Transformers’. arXiv, Aug. 16, 2021. doi: 10.48550/arXiv.2104.02057.
- [27] A. Dosovitskiy *et al.*, ‘An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale’. arXiv, Jun. 03, 2021. doi: 10.48550/arXiv.2010.11929.
- [28] T. Ridnik, E. Ben-Baruch, A. Noy, and L. Zelnik-Manor, ‘ImageNet-21K Pretraining for the Masses’. arXiv, Aug. 05, 2021. doi: 10.48550/arXiv.2104.10972.
- [29] C. Schuhmann *et al.*, ‘LAION-5B: An open large-scale dataset for training next generation image-text models’. arXiv, Oct. 15, 2022. doi: 10.48550/arXiv.2210.08402.
- [30] ‘Python Release Python 3.7.5’, Python.org. Accessed: Mar. 19, 2024. [Online]. Available: <https://www.python.org/downloads/release/python-375/>
- [31] R. Wightman, ‘PyTorch Image Models’, *GitHub repository*. GitHub, 2019. doi: 10.5281/zenodo.4414861.
- [32] L. F. Nakayama *et al.*, ‘A Brazilian Multilabel Ophthalmological Dataset (BRSET)’. PhysioNet. doi: 10.13026/XCXW-8198.
- [33] P. Porwal *et al.*, ‘Indian Diabetic Retinopathy Image Dataset (IDRiD): A Database for Diabetic Retinopathy Screening Research’, *Data*, vol. 3, no. 3, Art. no. 3, Sep. 2018, doi: 10.3390/data3030025.
- [34] M. Oquab *et al.*, ‘DINOv2: Learning Robust Visual Features without Supervision’. arXiv, Feb. 02, 2024. doi: 10.48550/arXiv.2304.07193.
- [35] ‘How to support the transition to AI-powered healthcare’, *Nat Med*, vol. 30, no. 3, pp. 609–610, Mar. 2024, doi: 10.1038/s41591-024-02897-9.
- [36] T. Darcet, M. Oquab, J. Mairal, and P. Bojanowski, ‘Vision Transformers Need Registers’. arXiv, Sep. 28, 2023. doi: 10.48550/arXiv.2309.16588.
- [37] C. Yang, J. Xu, S. De Mello, E. J. Crowley, and X. Wang, ‘GPViT: A High Resolution Non-Hierarchical Vision Transformer with Group Propagation’, arXiv.org. Accessed: Feb. 20, 2024. [Online]. Available: <https://arxiv.org/abs/2212.06795v2>
- [38] C. Ryali *et al.*, ‘Hiera: A Hierarchical Vision Transformer without the Bells-and-Whistles’, arXiv.org. Accessed: Feb. 20, 2024. [Online]. Available: <https://arxiv.org/abs/2306.00989v1>
- [39] C. de Vente *et al.*, ‘AIROGS: Artificial Intelligence for ROBust Glaucoma Screening Challenge’. arXiv, Feb. 10, 2023. doi: 10.48550/arXiv.2302.01738.

- [40] T. Li, Y. Gao, K. Wang, S. Guo, H. Liu, and H. Kang, 'Diagnostic assessment of deep learning algorithms for diabetic retinopathy screening', *Information Sciences*, vol. 501, pp. 511–522, Oct. 2019, doi: 10.1016/j.ins.2019.06.011.
- [41] I. Loshchilov and F. Hutter, 'Decoupled Weight Decay Regularization'. arXiv, Jan. 04, 2019. doi: 10.48550/arXiv.1711.05101.
- [42] I. Loshchilov and F. Hutter, 'SGDR: Stochastic Gradient Descent with Warm Restarts'. arXiv, May 03, 2017. doi: 10.48550/arXiv.1608.03983.
- [43] F. Pedregosa et al., 'Scikit-learn: Machine Learning in Python', *Journal of Machine Learning Research*, vol. 12, no. 85, pp. 2825–2830, 2011.
- [44] 'GPU Benchmarks for Deep Learning | Lambda'. Accessed: Feb. 14, 2024. [Online]. Available: <https://lambdalabs.com/gpu-benchmarks>
- [45] 'Record renewable energy output'. Accessed: Mar. 19, 2024. [Online]. Available: <http://www.gov.scot/news/record-renewable-energy-output/>
- [46] M. B. A. McDermott, L. H. Hansen, H. Zhang, G. Angelotti, and J. Gallifant, 'A Closer Look at AUROC and AUPRC under Class Imbalance'. arXiv, Jan. 11, 2024. doi: 10.48550/arXiv.2401.06091.
- [47] F. Wilcoxon, 'Individual Comparisons by Ranking Methods', *Biometrics Bulletin*, vol. 1, no. 6, pp. 80–83, 1945, doi: 10.2307/3001968.
- [48] J. Demšar, 'Statistical Comparisons of Classifiers over Multiple Data Sets', *J. Mach. Learn. Res.*, vol. 7, pp. 1–30, Dec. 2006.
- [49] A. J. Bishara and J. B. Hittner, 'Confidence intervals for correlations when data are not normal', *Behav Res*, vol. 49, no. 1, pp. 294–309, Feb. 2017, doi: 10.3758/s13428-016-0702-8.