

Interpretable Fully Convolutional Classification of Intrapapillary Capillary Loops for Real-Time Detection of Early Squamous Neoplasia

Luis C. Garcia-Peraza-Herrera¹ Martin Everson² Wenqi Li¹ Inmanol Luengo¹
Lorenz Berger⁵ Omer Ahmad¹ Laurence Lovat² Hsiu-Po Wang³ Wen-Lun
Wang⁴ Rehan Haidry² Danail Stoyanov¹ Tom Vercauteren¹ Sebastien
Ourselin¹

¹ Wellcome / EPSRC Centre for Interventional and Surgical Sciences, London, UK

² University College London Hospitals, London, UK

³ Gastroenterology National Taiwan University, Taipei, Taiwan

⁴ E-Da Cancer Hospital, Kaohsiung, Taiwan

⁵ Innersight Labs, London, UK

Abstract. In this work, we have concentrated our efforts on the interpretability of classification results coming from a fully convolutional neural network. Motivated by the classification of oesophageal tissue for real-time detection of early squamous neoplasia, the most frequent kind of oesophageal cancer in Asia, we present a new dataset and a novel deep learning method that by means of deep supervision and a newly introduced concept, the embedded Class Activation Map (eCAM), focuses on the interpretability of results as a design constraint of a convolutional network. We present a new approach to visualise attention that aims to give some insights on those areas of the oesophageal tissue that lead a network to conclude that the images belong to a particular class and compare them with those visual features employed by clinicians to produce a clinical diagnosis. In comparison to a baseline method which does not feature deep supervision but provides attention by grafting Class Activation Maps, we improve the F1-score from 87.3% to 92.7% and provide more detailed attention maps.

1 Introduction

Motivated by the clinical problem of intrapapillary capillary loops (IPCL) classification, we introduce a novel dataset containing 7046 frames from 17 patients (see table 1 for more details), and a novel unified framework for automatic feature extraction, classification and visual interpretability of results. We present a novel convolutional network architecture that focuses on the interpretability of the results as a design constraint for the network and serves as a baseline for quantitative comparison of results with future methods. We compare the visual features highlighted in the heatmaps produced by the network with those that are clinically relevant to produce a diagnosis.

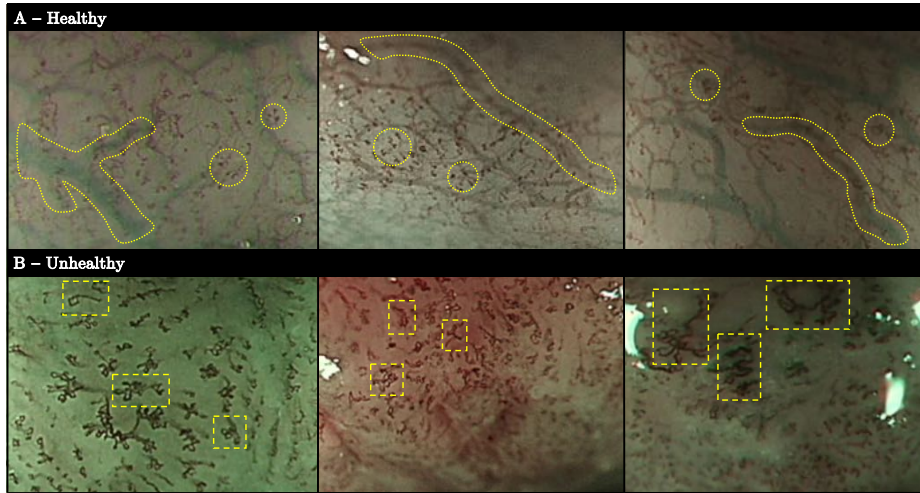


Fig. 1. Narrow-Band Imaging magnifying endoscopy of the oesophagus. (A) Frames extracted from surveillance endoscopies on several healthy subjects. In these patients mucosal vessels can be easily observed (dotted). In addition, interpapillary capillary loops (circles) are perceived as minuscule dots connected to an extremely thin filament. (B) Images from patients with abnormal interpapillary capillary loops suggesting carcinoma depth invasion. Microvessels are dilated and present unusual shape irregularities (rectangles).

In the Computer-Assisted Interventions (CAI) community labelled data is often scarce. Deep learning has become extremely popular due to its success in tasks such as classification and segmentation, but a large amount of data is typically required to capture the variability of the data across patients. As researchers in this area, we are also faced with additional challenges. Clinical collaborators are interested in interpreting the results coming from computer-assisted systems. That is, understanding the process followed by deep learning approaches to make a diagnosis. This includes analysing which features present in the images lead to a certain output and if those coincide with the ones that they analyse during clinical examination of the data. Conversely, it is also interesting to discover whether automatically extracted features are different from the ones currently used in clinical practice but can nonetheless lead to a correct diagnosis.

Using reduced datasets can potentially lead to models that do not generalise well. While it is true that there are efforts to build large scale CAI databases [1], in this paper, we have concentrated our efforts on the interpretability of classification results coming from a fully convolutional neural network trained on a small dataset.

Squamous cell carcinoma (SCC) is the most frequent kind of oesophageal cancer in Asia [2], presenting rapidly increasing numbers in the western world in recent decades too. Early diagnosis -and resection- play a key role to increase the chances of survival [3], as superficial lesions present low rates of lymphatic

dissemination. Detection is currently achieved by screening programs on high risk populations [2]. Narrow-Band Imaging magnifying endoscopy (NBI-ME) is the state-of-the-art technique employed for screening [4]. In addition to early diagnosis, a precise estimation of depth of invasion is crucial. Lesions that are closer to the oesophageal surface (mucosal layer) can be treated by minimally invasive endoscopic therapy rather than surgery [5].

NBI-ME facilitates the visualisation of micro-vascular patterns, called intra-papillary capillary loops (IPCL), which are linked to early squamous cell carcinoma and present focal, subtle, and easily missed visual features, particularly in centres with a low amount of cases. It has been also shown that the thickness and tortuosity of IPCL patterns is highly correlated with histological state and depth of invasion [4]. Hence, having an automated red-flag system that analyses each video frame in real-time could potentially help detect subtle IPCL patterns that might be difficult to distinguish by unspecialised endoscopists (see figure 1).

Recent work has explored different approaches to analyse the implicit attention mechanisms of convolutional neural networks. In [6], authors produce attention heatmaps as a linear combination of feature maps from the last convolutional layer. The weighting coefficients are extracted from the fully connected output neurons. Zhou et al. [7] allow for a fully convolutional classification by means of Global Average Pooling (GAP). When GAP is omitted (at inference), instead of a vector of class probabilities, a Class Activation Map (CAM) is automatically generated. In addition, these maps enable for accurate object localization, a task for which the network has not been trained for.

2 Materials and Methods

2.1 Deeply Supervised Embedded Class Activation Maps (eCAM)

There are various reasons why a fully convolutional classification is convenient for the proposed pretext task. Different endoscope processors provide images of varying resolutions. We aim for a flexible method that can generalise to different input sizes to simplify the preprocessing of data. It is also of interest to give the method versatility to process both full images or cropped patches. Furthermore, there are images of the oesophagus that can present unhealthy IPCL patterns only in certain areas. Hence, we seek for a method that could potentially have the ability to classify seamlessly both images and patches. GAP [7] has been shown as a feasible way to reduce the feature maps to a single value and still maintain a state-of-the-art classification accuracy.

As interpretability of the results is of utmost relevance, in addition to the classification score, we aim to obtain an attention heatmap that exposes to workings of the inference process and highlights those visual features that led the network conclude an image belongs to a certain class. This is relevant because it helps to check whether the network is paying attention to those parts of the image that clinical experts consider to be determinant to produce a valid diagnosis.

Furthermore, it serves as a validation mechanism that could point out possible problems in the learning process, for example, in case the network pays attention to areas of the image that are clinically irrelevant but happen to contain discriminative spurious visual features.

Class Activation Maps [7] are a recent attempt to produce meaningful attention heatmaps. They have shown to produce comparable to state-of-the-art localizations without re-training for the task. However, with this approach and the baseline network presented in figure 2, we achieve low resolution attention heatmaps that might allow to find a large object in scenes of daily life but lack the definition to illustrate attention in oesophageal NBI-ME images. In [8], the authors show that deep supervision is able to achieve superior results on segmentation of medical images. The reasons to opt for a similar approach are three-fold: fast convergence as gradients flow quicker to early layers of the network, ability to produce predictions at different resolutions, and improved quantitative classification results. Furthermore, we can take advantage of deep supervision to produce high resolution attention heatmaps (based on learnt deconvolutions that upsample the heatmaps to the original image size). As opposed to [7], we do not only aim to generate heatmaps that show the implicit attention of the network, but want to explicitly force the architecture to produce one attention map per class and use them to generate a classification prediction. We therefore introduce a new architecture that is fully convolutional, produces embedded attention heatmaps that allow for clinical interpretation of the results and works in real-time. The proposed method is shown in figure 2. The deep supervision mechanism is composed by several losses. As can be observed in figure 2, a loss

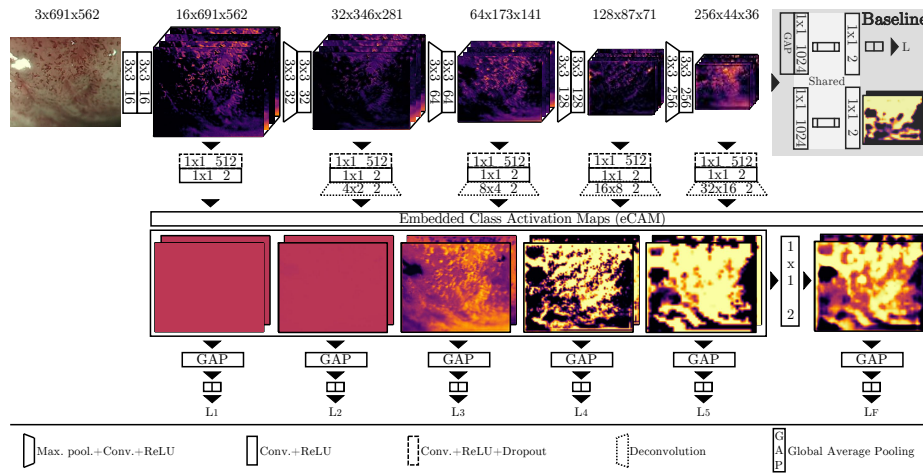


Fig. 2. Proposed convolutional network with multi-scale embedded Class Activation Maps. Baseline architecture (i.e. without deep supervision and classical CAM) shown shaded.

$L_s(\hat{\theta})$ is calculated using cross-entropy at each resolution scale:

$$L_s(\hat{\theta}) = - \sum_{k=1}^K g^{nk} \log \hat{p}_s^{nk} \quad (1)$$

where $\hat{\theta}$ are the network weights, $s \in \{1, \dots, 5\}$, g^{nk} is one when observation n belongs to class k and zero otherwise, and \hat{p}_s^{nk} is the predicted probability for observation n belonging to class k at scale s . $L_F(\hat{\theta})$ (see figure 2) has the same structure as $L_s(\hat{\theta})$ but the probability for each class comes from applying GAP to a weighted sum of the eCAM at all scales. That is, all the scale-dependent attention heatmaps are fused by means of a learnt 1x1 convolution so that the attention information at all scales is employed to produce a classification prediction. The final training loss to be minimised during the training process is

$$\mathcal{L}(\hat{\theta}) = \frac{1}{S+1} \left(L_F(\hat{\theta}) + \sum_{s=1}^S L_s(\hat{\theta}) \right) \quad (2)$$

where $S = 5$ scales.

2.2 Dataset

As we are introducing a new clinical problem, no dataset exists for the IPCL classification task. This novel dataset originates from 17 monocular videos (one video per patient) captured with two NBI-ME systems, OLYMPUS LUCERA CV-260 & CV-290 (Olympus Corporation, Tokyo, Japan). These videos have been recorded during routine screenings, and depict oesophageal recordings, starting on the stomach pit and finishing on the upper oesophageal sphincter. The oesophagus is cleaned prior to examination to expose clearly the mucosa.

The video sequences are cut to extract the useful parts of the procedure and sampled at 30fps as the endoscopists tend to perform rapid movements with the camera. The extracted frames are then quality controlled by an expert to discard those images that do not allow to perform a diagnosis with confidence. The final images are cropped so that no black corners or borders are left. All the pixels belong to oesophageal tissue. The labels have been matched with histological results from biopsies performed during the screening. The dataset contains 7046 frames, whose resolution ranges from 458x308 to 696x308 pixels. The dataset has been divided in three subsets, training, validation and testing. Each subset contains frames from different patients. To perform a thorough evaluation, five cross-validation folds have been created. Each fold contains a different draw of patients (see table 1 for more details).

All procedures performed in studies involving human participants were in accordance with the ethical standards of the local institutions and with the 1964 Helsinki declaration and its later amendments or comparable ethical standards.

Table 1. Number of frames for each cross-validation split. For each fold, the frames in the training, validation and testing sets belong to different patients. The letters (A) and (B) indicate whether the frames belong to class A -healthy- or B -unhealthy-.

Fold No.	Train. (A)	Val. (A)	Test. (A)	Train. (B)	Val. (B)	Val. (B)
1	2620	201	577	2803	258	587
2	1792	891	715	2205	739	704
3	1822	685	891	1549	1360	739
4	1792	715	891	1912	961	775
5	1559	685	1154	1754	743	1151
Average	1917	635	646	2045	812	791

2.3 Implementation details

Stochastic Gradient Descent (SGD) was the optimizer of choice. The training was performed with a fixed learning rate across training of $1e - 6$, momentum of 0.9, and a batch size of 1. A different CNN is trained for each dataset fold. All the networks are trained with a maximum number of iterations of $4 \times$ the number of images in the fold’s training set. Training weights are saved every 200 iterations and the best performing snapshot in validation set is selected for testing. CAFFE 1.0.0-rc5 [9] with CUDNN 5.1.10, CUDA 8.0.61, and NVIDIA driver 384.111 was the deep learning setup for development. The experiments were run on an Intel Core i7-4790K CPU @ 4.00GHz and an NVIDIA GeForce TITAN X (Pascal).

3 Results and Discussion

The proposed method achieves an average sensitivity and F1-score across dataset folds of 89.7% and 92.7% respectively in comparison with the baseline sensitivity and F1-score of 82.7% and 87.3% (see detailed quantitative evaluation on table 2). As shown in [10] deep supervision boosts accuracy by forcing the network to learn discriminant features at different resolutions. This is particularly relevant for endoscopy as features are visible or not depending on the distance from the camera to the oesophageal wall and the network has to be able to learn not only which features are useful at each scale but also how to fuse the predictions from different resolutions to achieve a correct classification. The prediction time interval ranges from 26.17ms for the smallest images in the dataset to 37.48ms for the largest ones.

In figure 3 we show different video frames and their corresponding eCAM. Only those from resolution levels L3, L4 and L5 and the multi-scale fused version are shown. The eCAM of resolution scale L1 and L2 are highly uncertain, as can be observed in figure 2. The reason possibly being two-fold. First, it is too early in the network and there are not enough filters (design constraint to achieve real-time) to capture the complexity of the disease. Second, receptive field being too small to capture discriminative visual features at those resolutions.

Table 2. Testing set classification results for the unhealthy class. Sensitivity, specificity, accuracy, precision and F1-score are reported.

Fold No.	Baseline					Proposed				
	Sens.	Spec.	Acc.	Prec.	F1	Sens.	Spec.	Acc.	Prec.	F1
1	77.5	99.8	88.6	99.8	87.2	80.4	92.0	86.2	91.1	85.4
2	39.2	99.9	69.8	99.6	56.3	78.1	99.7	89.0	99.6	87.6
3	100.0	97.1	98.4	96.6	98.3	100.0	95.9	97.7	95.2	97.6
4	96.6	97.9	97.3	97.5	97.1	99.4	97.3	98.3	97.0	98.2
5	100.0	95.6	97.8	95.8	97.8	90.6	99.6	95.1	99.5	94.9
Average	82.7	98.0	90.4	97.9	87.3	89.7	96.9	93.3	96.5	92.7

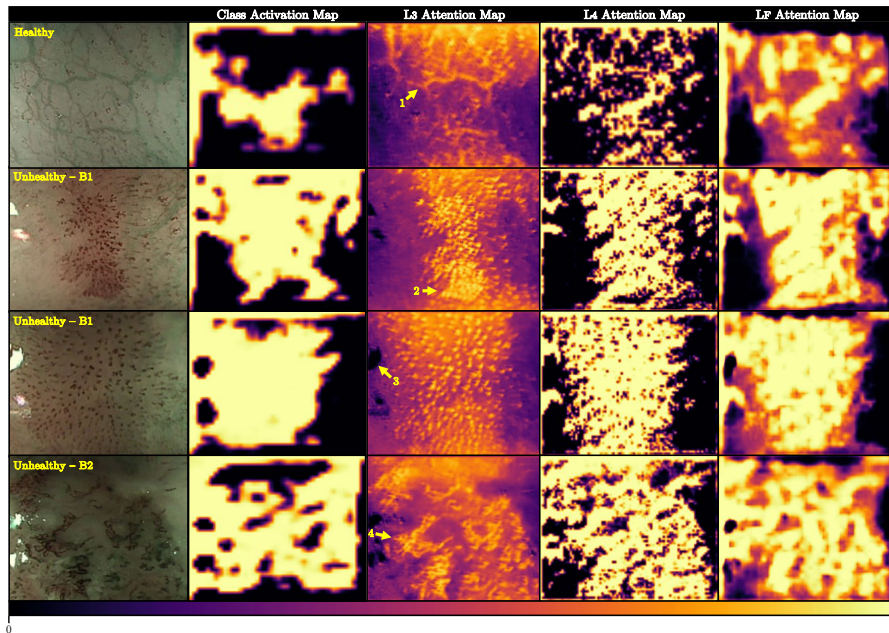


Fig. 3. Deeply supervised embedded Class Activation Maps. (1) Mucosal vessel erroneously highlighted in L_3 . (2) Densely populated IPCL area highlighted. (3) Specular reflection discarded. (4) *Star-shaped* irregular IPCL pattern highlighted.

Figure 3 shows several interesting visual features captured by eCAM. Specular reflections which are uninformative for tissue classification are discarded. The heatmaps are able to highlight both global and disperse unhealthy IPCL patterns and focalized areas of diseased tissue. Furthermore, despite recognising vessels as a matter of attention, the network does not seem to be able to discern that thick mucosal vessels are benign in healthy frames. Irregular *star-shaped* severe IPCL patterns (shown in the last row of figure 3) are also successfully highlighted as diseased. Large areas of healthy tissue are successfully recognised as benign as can be observed in the first row of figure 3.

4 Conclusion

Motivated by the problem of oesophageal IPCL pattern binary classification (healthy vs. unhealthy), we presented the first publicly available dataset for the task. We proposed a novel deeply supervised convolutional architecture that performs real-time fully convolutional classification achieving an average F1-score of 92.7%. We introduced the concept of embedded Class Activation Maps (eCAM) as a technique to force the network to capture and store visual attention maps and use them as source of information for classification. We showed that by means of deep supervision it is possible to obtain high quality heatmaps at the original resolution of the image. Future work will focus on the extension to multi-class detection of different unhealthy IPCL patterns.

References

1. Maier-Hein, L., Vedula, S.S., Speidel, S., Navab, N., Kikinis, R., Park, A., Eisenmann, M., Feussner, H., Forestier, G., Giannarou, S., Hashizume, M., Katic, D., Kenngott, H., Kranzfelder, M., Malpani, A., März, K., Neumuth, T., Padoy, N., Pugh, C., Schoch, N., Stoyanov, D., Taylor, R., Wagner, M., Hager, G.D., Jannin, P.: Surgical data science for next-generation interventions. *Nature Biomedical Engineering* **1**(9) (2017) 691–696
2. Wang, G.Q., Jiao, G.G., Chang, F.B., Fang, W.H., Song, J.X., Lu, N., Lin, D.M., Xie, Y.Q., Yang, L.: Long-term results of operation for 420 patients with early squamous cell esophageal carcinoma discovered by screening. *The Annals of Thoracic Surgery* **77**(5) (2004) 1740–1744
3. Endo, M., Kawano, T.: Detection and classification of early squamous cell esophageal cancer. *Diseases of the Esophagus* **10**(3) (1997) 155–158
4. Oyama, T., Inoue, H., Arima, M., Momma, K., Omori, T., Ishihara, R., Hirasawa, D., Takeuchi, M., Tomori, A., Goda, K.: Prediction of the invasion depth of superficial squamous cell carcinoma based on microvessel morphology: magnifying endoscopic classification of the Japan Esophageal Society. *Esophagus* **14**(2) (2017) 105–112
5. Ono, H.: Early gastric cancer: diagnosis, pathology, treatment techniques and treatment outcomes. *European Journal of Gastroenterology & Hepatology* (2006)
6. Cruz-Roa, A.A., Arevalo Ovalle, J.E., Madabhushi, A., González Osorio, F.A.: A Deep Learning Architecture for Image Representation, Visual Interpretability and Automated Basal-Cell Carcinoma Cancer Detection. In: MICCAI. (2013) 403–410
7. Zhou, B., Khosla, A., Lapedriza, A., Oliva, A., Torralba, A.: Learning Deep Features for Discriminative Localization. In: CVPR. (2016)
8. Garcia-Peraza-Herrera, L.C., Li, W., Fidon, L., Gruijthuijsen, C., Devreker, A., Attilakos, G., Deprest, J., Poorten, E.V., Stoyanov, D., Vercauteren, T., Ourselin, S.: ToolNet: Holistically-nested real-time segmentation of robotic surgical tools. In: IROS. (2017)
9. Jia, Y., Shelhamer, E., Donahue, J., Karayev, S., Long, J., Girshick, R., Guadarrama, S., Darrell, T.: Caffe: Convolutional Architecture for Fast Feature Embedding. In: Proceedings of the ACM International Conference on Multimedia. (2014) 675–678
10. Xie, S., Tu, Z.: Holistically-Nested Edge Detection. In: ICCV. (2015)