# Microscopic fine-grained instance classification through deep attention

Mengran Fan[1] Tapabrata Chakraborti[1] Eric I-Chao Chang[2] Yan Xu[2,3] and
Jens Rittscher[1]

[1] Institute of Biomedical Engineering, Dept. of Engg. Science, Univ. of Oxford, UK
`mengran.fan; tapabrata.chakraborty; jens.rittscher  @eng.ox.ac.uk`
[2] Microsoft Research, Beijing, China
`echang@microsoft.com`
[3] Department of Biology and Medicine, Beihang University, Beijing, China
`xuyan04@gmail.com`

**Abstract.** Fine-grained classification of microscopic image data with
limited samples is an open problem in computer vision and biomedi-
cal imaging. Deep learning based vision systems mostly deal with high
number of low-resolution images, whereas subtle detail in biomedical im-
ages require higher resolution. To bridge this gap, we propose a simple
yet effective deep network that performs two tasks simultaneously in an
end-to-end manner. First, it utilises a gated attention module that can
focus on multiple key instances at high resolution without extra anno-
tations or region proposals. Second, the global structural features and
local instance features are fused for final image level classification. The
result is a robust but lightweight end-to-end trainable deep network that
yields state-of-the-art results in two separate fine-grained multi-instance
biomedical image classification tasks: a benchmark breast cancer histol-
ogy dataset and our new fungi species mycology dataset. In addition,
we demonstrate the interpretability of the proposed model by visualising
the concordance of the learned features with clinically relevant features.

**Keywords:** Medical image classification · Deep attention mechanism

## 1 Introduction

Fine-grained image classification, which focuses on distinguishing subtle visual
differences between classes, is an open problem in biomedical image analysis.
Deep learning has led to a remarkable progress in fine-grained classification on
large-scale natural images [23,21,14]. Despite the important advances in com-
puter vision, it is usually challenging to achieve the same success on specific
biomedical image classification tasks [22,20]. To sum up, current methods mainly
face three challenges. Due to the cost of data acquisition and the limited avail-
ability of specimens, well-organised medical datasets in medical usually tend to
be small, which limits the representation ability of deep networks. The main rea-
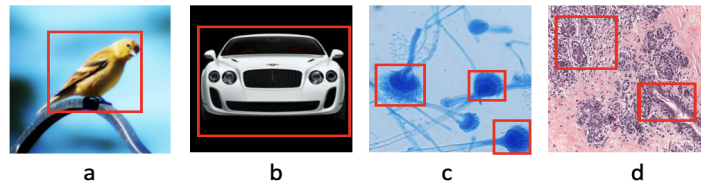son is that the current state-of-the-art convolutional neural networks (CNN) are

**Fig. 1: Classification challenges in biomedical imaging.** Compared to large-scale datasets on natural images (a-b), well-organised biomedical datasets (c-d) tend to be small and require expensive human expert annotations. Secondly, while we typically expect one centre-aligned instance (red box) in natural images, there are often multiple discriminative instances (red boxes) in biomedical images, which poses another challenge for feature learning. Finally, interpretability needs to be considered when developing a reliable medical image analysis system.

capable of extracting semantically meaningful features on large-scale datasets. When training data is limited deep networks may overfit and may bias the classification result on confounding background clutter.

Especially when working on the microscopic scale, multiple instances (e.g. glands, vessels, or crypts) need to be considered. This seriously restricts the adaptation of existing methods in the fine-grained classification of natural images. For instance, we randomly select samples (Fig.1 (a) - (b)) from the most popular fine-grained datasets in computer vision (CUB-200-2001 [16] and Stanford Cars [11]), where there is mostly one centre-aligned instance in an image. Although a large number of strategies have been proposed to detect the discriminative parts (e.g., head, belly for birds) in such images, the size and layout of the detected components are almost identical for each image. In comparison to natural images, biomedical images (Fig.1 (c) - (d)) may have a wide variety of discriminative instances (regions) with different sizes and densities, leading to more complicated structural information and a larger within-class variation. This motivates the need for investigating methods for building comprehensive and discriminative feature representations that can be applied in this domain. Thirdly, apart from the accurate prediction, the interpretability also plays a crucial role in a reliable medical image classification system [19]. In this work, we propose a novel attention-based classification network that is capable of jointly localising discriminative instances and enhancing consistent fine-grained feature learning in an end-to-end fashion. The main contributions of this paper are: (1) A lightweight gated attention module where the most discriminative instances can be localised simultaneously without requiring any part annotations or redundant region proposals. (2) A multi-task learning scheme that dynamically controls the weights of member modules and enforces the network to learn consistent instance-level features. (3) Improved interpretability of learned features when compared with features used by human experts for decision making.
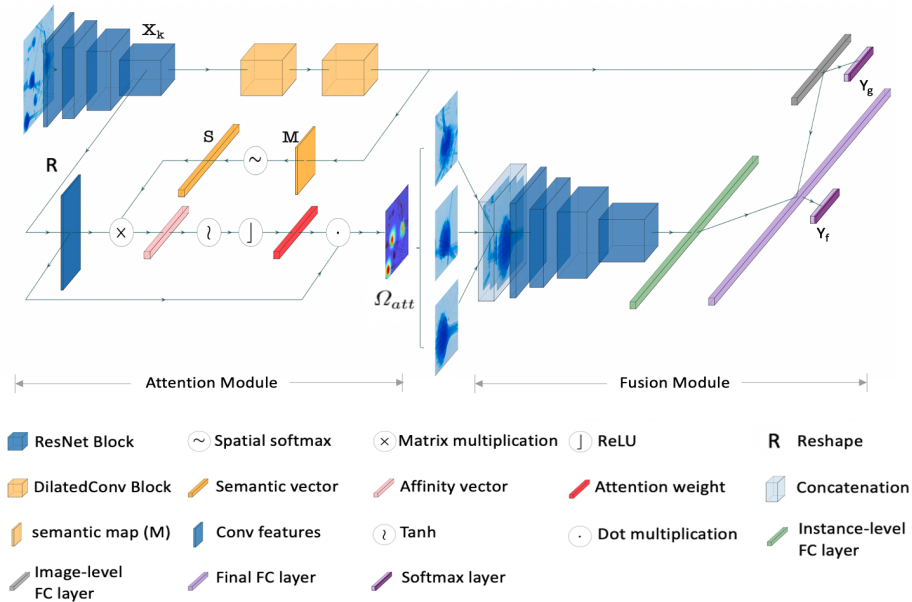
**Fig. 2: Framework for the proposed multiple instance fine-grained classification pipeline.** The proposed network consists of two main modules: the attention module and the feature fusion module. The details of the attention module and definition of the variables are explained in Section 2.1. In the fusion module, we threshold the attention map to generate a binary mask, crop these patches from the input image, resize them to a predefined size and feed them into a shared CNN model. Finally, we concatenate all instance-level features with image-level features for final prediction. The predictions in the inference stage are also conducted in this way.

## 2    Methodology

We propose a novel fine-grained multi-instance classification scheme (Fig.2) that consists of two main modules: (i) a novel gated attention mechanism for discriminative instances localisation; and (ii) a feature fusion strategy that consolidates the global and local features to facilitate the final predictions.

### 2.1    Gated attention module

In natural images, noisy background or irrelevant regions are highly variable and can be therefore naturally discarded by deep neural networks. However, for medical datasets with limited training samples, this is more difficult to achieve. Here, we propose a gated attention mechanism that is aimed to filter out the confounding channels and localise the most discriminative instances without extra part supervision or redundant region proposals. As shown in Fig.2, the module first summarises a $2D$ semantic map $S \in \mathbb{R}^{H \times W}$ from the input convolutional feature maps $X \in \mathbb{R}^{C \times H \times W}$. Furthermore, we use the high-level semantic information

to filter out the confounding channels based on the pairwise correspondences between each input channel and the generated semantic map, thus suppressing the irrelevant background and localising the most discriminative instances.

**Semantic Modelling.** Since spatial attention weights computed for each query position are almost the same for different tasks [3], we extract a global spatial attention map from the input features, which are shared by all query positions within an image. For multi-level semantics understanding, we first apply two dilated convolution blocks [4] to the output of the feature extractor $X_k$. The set of multi-scale features are compressed by computing the sum of all channels $M = \sum_{k=1}^{C} W^T X_k$ where $W^T$ is the weights of dilated blocks. This channel compression rests on the assumption that if the region is activated on most channels, the region tends to be more discriminative and to have higher likelihood of being part of the object of interest. The final semantic map $S$ is generated by applying a spatial softmax layer that performs the softmax operation over all feature points in the aggregated map $M$, resulting in a probability distribution that roughly indicates the regions of the most discriminative instances:

$$S_{i,j} = \frac{\exp\left(M_{i,j}\right)}{\sum_{l=1}^{H} \sum_{k=1}^{W} \exp\left(M(l,k)\right)} \tag{1}$$

**Gated Mechanism.** To measure the discriminability of each channel, we capture the spatial correspondence scores via conducting matrix multiplication over $X^T$ and $S$, where $X^T$ is the original convolutional feature maps with the shape of $c \times hw$. For example, $X_k$ is the $k^{th}$ channel of the input feature maps, containing its specific semantic responses. So $X_k^T S$ is the importance coefficient that indicates the semantic representation power of this channel. Our method is different from traditional channel-based self-attention mechanisms [23,6] that usually directly capture the pairwise inter-channel dependencies by calculating $X^T X$. In order to enhance the specific semantics, we summarise a global high-level semantic map and use it as a template to quantify the representation capability of each channel by $X^T S$. Therefore, we apply such a mechanism to obtain the $1D$ coefficient vector $X^T S \in \mathbb{R}^{C \times 1}$ rather than $X^T X \in \mathbb{R}^{C \times C}$.

To filter out the confounding channels, the hyperbolic tangent ($tanh$) and ReLU activation functions are used to normalise the discriminability coefficient among all channels. As a result, a set of gated weights is obtained, selects channels that look at the most discriminative regions. In particular, the gated weight is approximately 1 for the most informative channels, and approximately 0 for the channels highlighting the irrelevant background (Fig. 3). The gated activation layer can also be regarded as a filter which enforces the model to ignore the confounding channels and pay attention to more informative channels. Consequently, to let the attention module focus on multiple instances, we model the final attention map as a gated average of the outputs of the original channels.

$$\Omega_{att} = \frac{1}{C} \sum_{k=1}^{C} (X_k^T \odot ReLU(\tanh\left(X_k^T \otimes S\right))) \tag{2}$$

### 2.2   Multi-task Loss Function

Different from traditional two-stage frameworks consisting of two separate networks, the multi-task loss aims to enable the model to jointly learn multi-instance localisation and image classification in an end-to-end fashion. Specifically, our network is optimised by a global attention loss and a final fusion loss:

$$\mathcal{L} = \lambda L_{\mathcal{G}}(Y_g, Y^\star) + (1 - \lambda)(L_{\mathcal{F}}(Y_f, Y^\star)) \tag{3}$$

where $L_{\mathcal{G}}$ and $L_{\mathcal{F}}$ are standard cross entropy losses with respect to the outputs of the global image-level network $Y_g$ and the proposed multi-instance fusion network $Y_f$, respectively. $Y^\star$ represents the ground truth label and the parameter $\lambda$ is initialised as 1 and gradually decreased during training. As a result, the network initially focuses on extracting global image-level features, and increases the contribution of discriminative instance-level features during training.

## 3   Evaluation

All input images were resized to $224 \times 224$, and a Resnet-18 was used to extract global image-level information from down-sampled images. After instance localisation, extracted patches were scaled to $336 \times 336$, and fed to a Resnet-50 for final image-level prediction. Other CNN architectures could be used instead. To improve training efficiency, pre-trained weights from the ImageNet dataset were used for initialisation. Mini batch size was set to 16. We used the stochastic gradient descent (SGD) optimiser with an initial learning rate of 0.05 that was multiplied by 0.1 after every 50 epochs. The initial weight score $\lambda$ in the loss function is 1, and reduced by 0.1 after every 20 training epochs. The publicly available MXNet library was used to implemented the model, training was performed on two NVIDIA GeForece 1080 Ti GPUs.

**Evaluation and performance analysis on new fungi species dataset.**To the best of our knowledge, this is the first attempt for bringing deep learning based approaches to fungal species identification. 2151 microscopy images from 59 patients were collected in collaboration with the Peking Union Medical College Hospital. In this dataset, we particularly focus on five most common species involved in human disease: (1)Aspergillus fumigatus, (2)Aspergillus flavus, (3)Aspergillus niger, (4)Aspergillus terreus and (5)Aspergillus nidulans. We provide quantitative results and compare it with recent competing methods. We also benchmark the performance of the novel gated attention mechanism with other attention schemes. For all experiments, we randomly split the samples in each class in a ratio of $1 : 3$ for constructing testing and training sets.

**Quantitative comparison with competing methods** We evaluated the effectiveness of the proposed method by comparing it with several state-of-the-art fine-grained classification methods. All of the compared methods were trained with the same backbone network and computing environment. From the comparison shown in Table 1, we observe that our method achieves the best performance when compared with other fine-grained classification methods.

**Table 1: Results on Fungi species dataset.**

| Methods | Accuracy |
|---------|----------|
| Resnet-50 [8] | 0.907 |
| Residual attention [17] | 0.867 |
| Attend & Rectify[14] | 0.871 |
| Trilinear attention [23] | 0.883 |
| NTS Network [21] | 0.914 |
| **Our method** | **0.951** |

**Table 2: Comparison of attention mechanisms.**

| Attention Mechanisms | Accuracy |
|----------------------|----------|
| Spatial attention [18] | 0.859 |
| Channel-wise attention [23] | 0.883 |
| Dual attention [6] | 0.901 |
| Squeeze-Excitation attention [9] | 0.939 |
| Global Context attention [3] | 0.937 |
| **Our Gated Attention** | **0.951** |

**Evaluation of gated attention mechanism.** To measure the effectiveness of our gated attention module we compared it with other existing attention mechanisms but using the same sampling strategy, feature fusion scheme and loss function. We only modified the attention-based instance localisation module in the baseline model, and investigate the performance of different attention mechanisms. Table 2 shows the results of integrating different attention modules in our classification framework, and Fig. 3 depicts a visualisation example of each step in our gated attention module. Our gated attention mechanism not only outperforms all other attention modules but it also suppresses the confounding information, demonstrating the effectiveness and localisation ability.

**Table 3: Results on breast cancer dataset.**

| Methods | Accuracy | Methods | Accuracy |
|---------|----------|---------|----------|
| Vgg19 [10] | 0.925 | Inception-v3 [10] | 0.913 |
| DenseNet-161 [10] | 0.940 | Model Fusion [13] | 0.925 |
| AlexNet [12] | 0.813 | ResNet-152 [2] | 0.830 |
| RFSVM-All [2] | 0.930 | Ensemble [15] | 0.825 |
| Refined Ensemble [15] | 0.875 | Two-stage network [7] | 0.850 |
| Hybrid deep network [20] | 0.913 | **Our method** | **0.970** |

**Evaluation on Breast Cancer Histology images.** The BreAst Cancer Histology images (BACH) benchmark dataset [1] is used to investigate the method's ability for histology images. This dataset consists of 400 high-resolution ( $2018 \times 1356$ ) Hematoxylin and eosin stained microscopy images, with an even distribution over four classes. Each image is labeled as one of four types: 1) normal, 2) benign, 3) in situ carcinoma and 4) invasive carcinoma, according to the predominant tissue type. We randomly perform a $75\% - 25\%$ split for training and testing. Table 3 shows the classification results on breast cancer histology images. We compared the best classification accuracy over several advanced methods in the case of the 400 images provided by the challenge organizer. Our approach achieves the best classification performance with 0.970, showing that
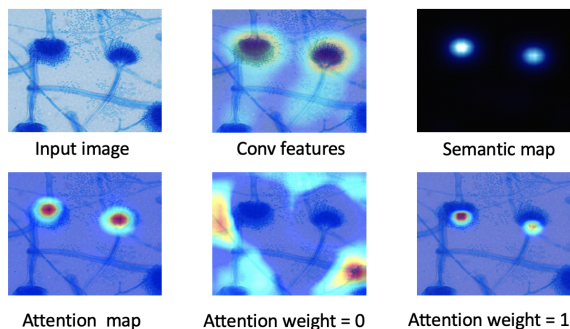
**Fig. 3: Visualization of the gated attention module for one sample image.** The top row shows the input image, original convolutional feature map and the corresponding semantic map (defined in Eq.1). The bottom row shows the final attention map and two representative channels with the lowest and highest attention.

our network can be effectively applied to the classification tasks of histology images.

**Interpretability and alignment with clinical background.** By analysing the concordance of the learned attention maps with well established visual clues used by human experts we evaluate their interpretability. The reader can easily appreciate the importance of this in addition to accuracy in results. **Fungal species:** In clinical practice, key criteria [24,5] are a range of morphological features associated with the structure of conidial heads, especially the colour, size and shape of vesicles (Fig.4 (b)). Fig.4 (a) shows the sample images and corresponding attention maps of each specie. Our attention maps consistently highlight the relevance of these vesicle patterns. **Breast cancer:** A normal healthy breast duct is made up of layers of inner epithelial cells, outer myoepithelial cells and a basement membrane (see Fig.5 (b)). In the case of *in situ* carcinomas, the proliferating cancer cells are restrained inside the basement membrane, whereas the cancer cells break out of the walls and invade the surrounding breast tissue in invasive cases. Thus, the intactness of the basement membrane is diagnostic relevance. To evaluate the effectiveness of learned features, bounding box annotations were generated by an expert pathologist on 100 test images. Overall, 72% of the bounding boxes are covered by our network and selected examples of *in situ* carcinomas are shown in Fig.5 (a).

## 4   Conclusion

We present a simple yet effective end-to-end deep architecture that addresses the problem of fine-grained multi-instance classification from biomedical images at high resolution. It achieves this by first using a lightweight gated attention mechanism that detects multiple key instances and then combining the global structure and local instance features for a final image level classification. The
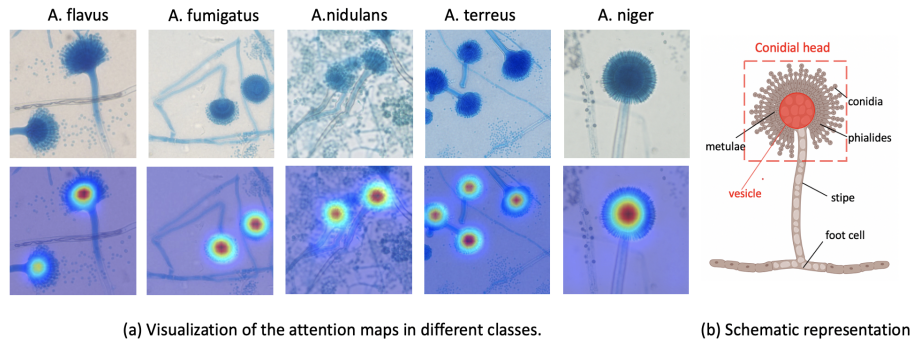
(a) Visualization of the attention maps in different classes.                    (b) Schematic representation

**Fig. 4: Clinical alignment on fungi dataset.** Clinicians mainly rely on the morphology assessment of conidial heads, especially vesicles (shown in (b)), to differentiate Aspergillus species. The attention maps (shown in (a)) generated by the proposed network consistently match the guideline for clinical decision making.
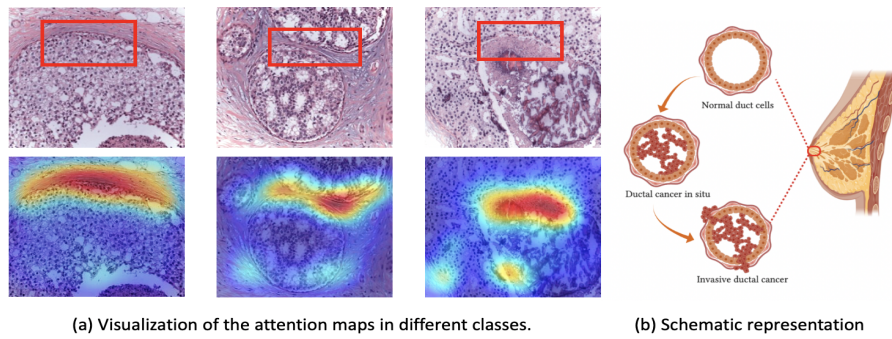


(a) Visualization of the attention maps in different classes.                    (b) Schematic representation

**Fig. 5: Clinical alignment on breast cancer dataset.** The first row in (a) shows the examples of *in situ* carcinomas with bounding box annotations. The attention maps shown in the second row consistently focus on the membrane boundaries, covering the human annotations .

proposed network is evaluated on a new fungi species classification dataset and a publicly available breast cancer dataset and achieves state-of-the-art performance. We also demonstrate in details the scope of our method as an interpretable model by showing the strong alignment of the learned features with well documented visual clues used by human subject matter experts.

# References

1. Aresta, G., Araújo, T., Kwok, S., Chennamsetty, S.S., Safwan, M., Alex, V., Marami, B., Prastawa, M., Chan, M., Donovan, M., et al.: Bach: Grand challenge on breast cancer histology images. Medical image analysis (2019)
2. Cao, H., Bernard, S., Heutte, L., Sabourin, R.: Improve the performance of transfer learning without fine-tuning using dissimilarity-based multi-view learning for breast cancer histology images. In: International Conference Image Analysis and Recognition. pp. 779–787. Springer (2018)
3. Cao, Y., Xu, J., Lin, S., Wei, F., Hu, H.: Gcnet: Non-local networks meet squeeze-excitation networks and beyond. arXiv preprint arXiv:1904.11492 (2019)
4. Chen, L.C., Papandreou, G., Kokkinos, I., Murphy, K., Yuille, A.L.: Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. IEEE transactions on pattern analysis and machine intelligence **40**(4), 834–848 (2017)
5. Diba, K., Kordbacheh, P., Mirhendi, S., Rezaie, S., Mahmoudi, M.: Identification of aspergillus species using morphological characteristics. Pakistan journal of medical sciences **23**(6),  867 (2007)
6. Fu, J., Liu, J., Tian, H., Li, Y., Bao, Y., Fang, Z., Lu, H.: Dual attention network for scene segmentation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 3146–3154 (2019)
7. Golatkar, A., Anand, D., Sethi, A.: Classification of breast cancer histology using deep learning. In: International Conference Image Analysis and Recognition. pp. 837–844. Springer (2018)
8. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 770–778 (2016)
9. Hu, J., Shen, L., Sun, G.: Squeeze-and-excitation networks. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 7132–7141 (2018)
10. Kohl, M., Walz, C., Ludwig, F., Braunewell, S., Baust, M.: Assessment of breast cancer histology using densely connected convolutional networks. In: International Conference Image Analysis and Recognition. pp. 903–913. Springer (2018)
11. Krause, J., Stark, M., Deng, J., Fei-Fei, L.: 3d object representations for fine-grained categorization. In: Proceedings of the IEEE international conference on computer vision workshops. pp. 554–561 (2013)
12. Nawaz, W., Ahmed, S., Tahir, A., Khan, H.A.: Classification of breast cancer histology images using alexnet. In: International Conference Image Analysis and Recognition. pp. 869–876. Springer (2018)
13. Rakhlin, A., Shvets, A., Iglovikov, V., Kalinin, A.A.: Deep convolutional neural networks for breast cancer histology image analysis. In: International Conference Image Analysis and Recognition. pp. 737–744. Springer (2018)
14. Rodríguez, P., Gonfaus, J.M., Cucurull, G., XavierRoca, F., Gonzalez, J.: Attend and rectify: a gated attention mechanism for fine-grained recovery. In: Proceedings of the European Conference on Computer Vision (ECCV). pp. 349–364 (2018)
15. Vang, Y.S., Chen, Z., Xie, X.: Deep learning framework for multi-class breast cancer histology image classification. In: International Conference Image Analysis and Recognition. pp. 914–922. Springer (2018)
16. Wah, C., Branson, S., Welinder, P., Perona, P., Belongie, S.: The caltech-ucsd birds-200-2011 dataset. Technical Report **CNS-TR-2011-201** (2011)

17. Wang, F., Jiang, M., Qian, C., Yang, S., Li, C., Zhang, H., Wang, X., Tang, X.: Residual attention network for image classification. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 3156–3164 (2017)
18. Wang, X., Girshick, R., Gupta, A., He, K.: Non-local neural networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 7794–7803 (2018)
19. Weese, J., Lorenz, C.: Four challenges in medical image analysis from an industrial perspective. Medical Image Analysis **33**, 44–49 (2016)
20. Yan, R., Ren, F., Wang, Z., Wang, L., Zhang, T., Liu, Y., Rao, X., Zheng, C., Zhang, F.: Breast cancer histopathological image classification using a hybrid deep neural network. Methods (2019)
21. Yang, Z., Luo, T., Wang, D., Hu, Z., Gao, J., Wang, L.: Learning to navigate for fine-grained classification. In: Proceedings of the European Conference on Computer Vision (ECCV). pp. 420–435 (2018)
22. Zhang, J., Xie, Y., Wu, Q., Xia, Y.: Medical image classification using synergic deep learning. Medical image analysis **54**, 10–19 (2019)
23. Zheng, H., Fu, J., Zha, Z.J., Luo, J.: Looking for the devil in the details: Learning trilinear attention sampling network for fine-grained image recognition. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 5012–5021 (2019)
24. Zulkifli, N.A., Zakaria, L.: Morphological and molecular diversity of aspergillus from corn grain used as livestock feed. HAYATI Journal of Biosciences **24**(1), 26–34 (2017)