

A Robust and Effective Approach Towards Accurate Metastasis Detection and pN-stage Classification in Breast Cancer

Byungjae Lee and Kyunghyun Paeng

Lunit inc., Seoul, South Korea
{jaylee,khpaeng}@lunit.io

Abstract. Predicting TNM stage is the major determinant of breast cancer prognosis and treatment. The essential part of TNM stage classification is whether the cancer has metastasized to the regional lymph nodes (N-stage). Pathologic N-stage (pN-stage) is commonly performed by pathologists detecting metastasis in histological slides. However, this diagnostic procedure is prone to misinterpretation and would normally require extensive time by pathologists because of the sheer volume of data that needs a thorough review. Automated detection of lymph node metastasis and pN-stage prediction has a great potential to reduce their workload and help the pathologist. Recent advances in convolutional neural networks (CNN) have shown significant improvements in histological slide analysis, but accuracy is not optimized because of the difficulty in the handling of gigapixel images. In this paper, we propose a robust method for metastasis detection and pN-stage classification in breast cancer from multiple gigapixel pathology images in an effective way. pN-stage is predicted by combining patch-level CNN based metastasis detector and slide-level lymph node classifier. The proposed framework achieves a state-of-the-art quadratic weighted kappa score of 0.9203 on the Camelyon17 dataset, outperforming the previous winning method of the Camelyon17 challenge.

Keywords: Camelyon17, Convolutional neural networks, Deep learning, Metastasis detection, pN-stage classification, Breast cancer

1 Introduction

When cancer is first diagnosed, the first and most important step is staging of the cancer by using the TNM staging system [1], the most commonly used system. Invasion to lymph nodes, highly predictive of recurrence [2], is evaluated by pathologists (pN-stage) via detection of tumor lesions in lymph node histology slides from a surgically resected tissue. This diagnostic procedure is prone to misinterpretation and would normally require extensive time by pathologists because of the sheer volume of data that needs a thorough review. Automated detection of lymph node metastasis and pN-stage prediction has the potential to significantly elevate the efficiency and diagnostic accuracy of pathologists for one of the most critical diagnostic process of breast cancer.

In the last few years, considerable improvements have been emerged in the computer vision task using CNN [3]. Followed by this paradigm, CNN based computer assisted metastasis detection has been proposed in recent years [4,5,6]. However, recent approaches metastasis detection in whole slide images have shown the difficulty in handling gigapixel images [4,5,6]. Furthermore, pN-stage classification requires handling multiple gigapixel images.

In this paper, we introduce a robust method to predict pathologic N-stage (pN-stage) from whole slide pathology images. For the robust performance, we effectively handle multiple gigapixel images in order to integrate CNN into pN-stage prediction framework such as balanced patch sampling, patch augmentation, stain color augmentation, 2-stage fine-tuning and overlap tiling strategy. We achieved patient-level quadratic weighted kappa score 0.9203 on the Camelyon17 test set which it yields the new state-of-the-art record on Camelyon17 leaderboard [7].

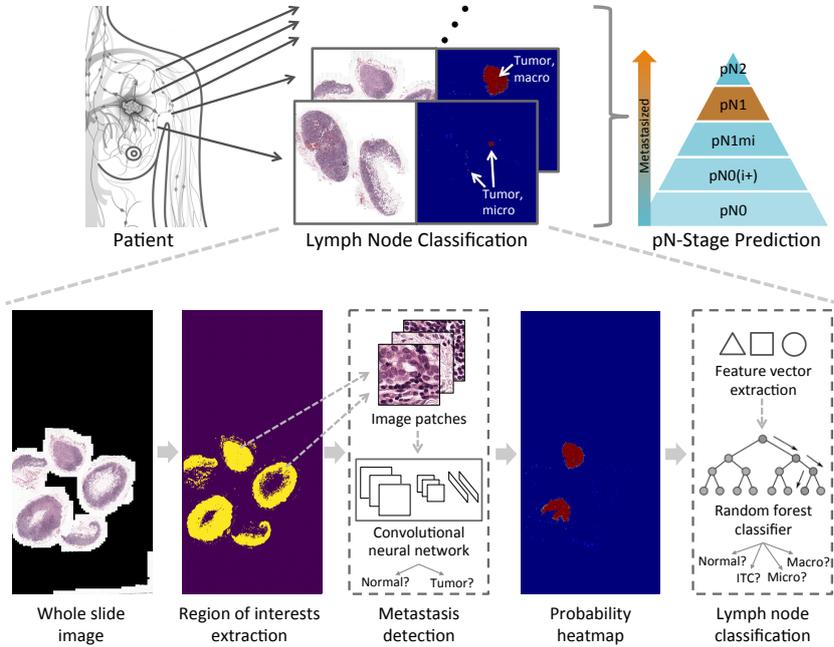


Fig. 1. Overall architecture of our pN-stage prediction framework.

2 Methodology

Fig. 1 shows the overall scheme of our proposed framework. First, ROI extraction module proposes candidate tissue regions from whole slide images. Second, CNN-based metastasis detection module predicts cancer metastasis within extracted ROIs. Third, the predicted scores extracted from ROI are converted to a feature vector based on the morphological and geometrical information which is used to

build a slide-level lymph node classifier. Patient-level pN-stage is determined by aggregating slide-level predictions with given rules [7].

2.1 Regions of Interests Extraction

A whole slide image (WSI) is approximately 200000×100000 pixels on the highest resolution level. Accurate tissue region extraction algorithms can save computation time and reduce false positives from noisy background area. In order to extract tissue regions from the WSIs, Otsu threshold [8] or gray value threshold is commonly used in recent studies [4,5,6]. We decide to use gray value threshold method which shows superior performance in our experiments.

2.2 Metastasis Detection

Some annotated metastasis regions include non-metastasis area since accurate pixel-level annotation is difficult in gigapixel WSIs [5]. We build a large scale dataset by extracting small patches from WSIs to deal with those noisy labels. After the ROIs are found from WSIs as described in Section 2.1, we extract 256×256 patches within ROIs with stride 128 pixels. We label a patch as tumor if over 75% pixels in the patch are annotated as a tumor. Our metastasis detection module is based on the well-known CNN architecture ResNet101 [3] for patch classification to discriminate between tumor and non-tumor patches.

Although the proposed method seems straightforward, we need to effectively handle gigapixel WSIs to integrate CNN into pN-stage prediction framework for the robust performance, as described below.

Balanced Patch Sampling The areas corresponding to tumor regions often covered only a minor proportion of the total slide area, contributing to a large patch-level imbalance. To deal with this imbalance, we followed similar patch sampling approach used in [5]. In detail, we sample the same number of tumor/normal patches where patches are sampled from each slide with uniform distribution.

Patch Augmentation There are only 400 WSIs in Camelyon16 dataset and 500 WSIs in Camelyon17 **train** set. Patches sampled from same WSI exhibit similar data property, which is prone to overfitting. We perform extensive data augmentation at the training step to overcome small number of WSIs. Since the classes of histopathology image exhibit rotational symmetry, we include patch augmentation by randomly rotating over angles between 0 and 360, and random left-right flipping. Details are shown in Table 1.

Table 1. Patch augmentation details.

Methods	Details
Translation	random x, y offset in $[-8, 8]$
Left/right flip	with 0.5 probability
Rotation	random angle in $[0, 360)$

Stain Color Augmentation To combat the variety of hematoxylin and eosin (H&E) stained color because of chemical preparation difference per slide, extensive color augmentation is performed by applying random hue, saturation, brightness, and contrast as described in Table 2. CNN model becomes robust against stain color variety by applying stain color augmentation at the training step.

Table 2. Stain color augmentation details.

Methods	Details
Hue	random delta in [-0.04, 0.04]
Saturation	random saturation factor in [0.75, 1.25]
Brightness	random delta in [-0.25, 0.25]
Contrast	random contrast factor in [0.25, 1.75]

2-Stage Fine-Tuning Camelyon16 and Camelyon17 dataset are collected from different medical centers. Each center may use different slide scanners, different scanning settings, difference tissue staining conditions. We handle this multi-center variation by applying the 2-stage fine-tuning strategy. First, we fine-tune CNN with the union set of Camelyon16 and Camelyon17 and then fine-tune CNN again with only Camelyon17 set. The fine-tuned model becomes robust against multi-center variation between Camelyon16 and Camelyon17 set.

Overlap Tiling Strategy In the prediction stage, probability heatmap is generated by the trained CNN based metastasis detector. A straightforward way to generate a heatmap from WSI is separating WSI into patch size tiles and merging patch level predictions from each tile. However, this simple strategy provides insufficient performance. Instead, we use similar overlap-tile strategy [9] for dense heatmap from tiled WSI. As shown in Fig. 2, the probability heatmap generated by overlap-tile strategy provides denser heatmap than straightforward tiling strategy even though the same classifier is used. By default, we used 50% overlapped tiles shown in Fig. 2(c).

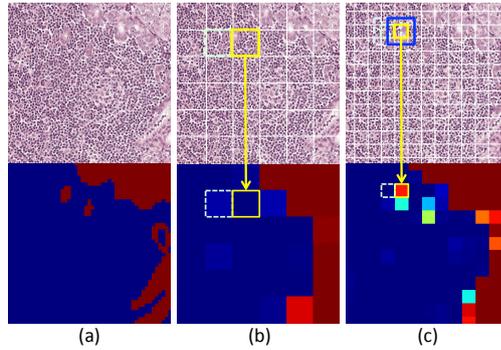


Fig. 2. Tiling strategy for dense heatmap. (a) A ground truth; (b) Straightforward tiling strategy; (c) Overlap-tile strategy.

As shown in Fig. 2, the probability heatmap generated by overlap-tile strategy provides denser heatmap than straightforward tiling strategy even though the same classifier is used. By default, we used 50% overlapped tiles shown in Fig. 2(c).

2.3 Lymph Node Classification

To determine each patient’s pN-stage, multiple lymph node slides should be classified into four classes (Normal, Isolated tumor cells (ITC), Micro, Macro). For each lymph node WSI, we obtain the $128 \times$ down-sampled tumor probability heatmap through the CNN based metastasis detector (Section 2.2). Each heatmap is converted into a feature vector which is used to build a slide level lymph node classifier. We define 11 types of features based on the morphological and geometrical information. By using converted features, random forest classifier [10]

is trained to automatically classify the lymph node into four classes. Finally, each patient’s pN-stage is determined by aggregating all lymph node predictions with the given rule [7]. We followed the Camelyon17’s simplified version of the pN-staging system (pN0, pN0(i+), pN1mi, pN1, pN2) [7].

3 Experiments

3.1 Dataset

We evaluate our framework on Camelyon16 [6] and Camelyon17 [7] dataset. The Camelyon16 dataset contains 400 WSIs with region annotations for all its metastasis slides. The Camelyon17 dataset contains 1000 WSIs with 5 slides per patient: 500 slides for the **train** set, 500 slides for the **test** set. The **train** set consists of the slide level metastasis annotation. There are 3 categories of lymph node metastasis: Macro (Metastases greater than 2.0 mm), Micro (metastasis greater than 0.2 mm or more than 200 cells, but smaller than 2.0 mm), and ITC (single tumor cells or a cluster of tumor cells smaller than 0.2mm or less than 200 cells).

Since the Camelyon17 set provides only 50 slides with lesion-level annotations in **train** set, we split 100 patients (total 500 WSIs since each patient provides 5 WSIs) into 43 patients for the Camelyon17 **train-M** set to train metastasis detection module, 57 patients for the Camelyon17 **train-L** set to train lymph node classification module. In detail, if patient’s any slide include lesion-level annotation, we allocate that patient as a Camelyon17 **train-M** set. Other patients are allocated as a Camelyon17 **train-L** set. As shown in Table 3, our split strategy separates similar data distribution between them in terms of the medical centers and metastasis types.

Table 3. Details of our Camelyon17 dataset split.

Dataset	# of patients per each pN-stage					Total
	pN0	pN0(i+)	pN1mi	pN1	pN2	
Camelyon17 train-M	0	9	11	14	9	43
Camelyon17 train-L	24	3	9	11	10	57
Dataset	# of patients per each medical center					Total
	Center1	Center2	Center3	Center4	Center5	
Camelyon17 train-M	7	8	9	10	9	43
Camelyon17 train-L	13	12	11	10	11	57
Dataset	# of WSIs per each metastasis type				Total	
	Negative	ITC	Micro	Macro		
Camelyon17 train-M	110	26	35	44	215	
Camelyon17 train-L	203	9	29	44	285	

3.2 Evaluation Metrics

Metastasis Detection Evaluation We used the Camelyon16 evaluation metric [6] on the Camelyon16 dataset to validate metastasis detection module performance. Camelyon16 evaluation metric consists of two metrics, the area under receiver operating characteristic (AUC) to evaluate the slide-level classification and the FROC to evaluate the lesion-level detection and localization.

pN-stage Classification Evaluation To evaluate pN-stage classification, we used the Camelyon17 evaluation metric [7], patient-level five-class quadratic weighted kappa where the classes are the pN-stages. Slide-level lymph node classification accuracy is also measured to validate lymph node classification module performance.

3.3 Experimental Details

ROI Extraction Module For the type of ROI extraction between Otsu threshold and gray value threshold, we determined to use gray value threshold method which is obtained a better performance on Camelyon16 **train** set. In detail, we convert RGB to gray from $32\times$ down-sampled WSI and then extract tissue regions by thresholding gray value > 0.8 .

Metastasis Detection Module During training and inference, we extracted 256×256 patches from WSIs at the highest magnification level of $0.243\ \mu\text{m}/\text{pixel}$ resolution. For training of the patch-level CNN based classifier, 400 WSIs from Camelyon16 dataset and 160 WSIs from Camelyon17 **train** set are used as shown in Table 4. Total 1,430K tumor patches and 43,700K normal patches are extracted.

We trained ResNet101 [3] with initial parameters from ImageNet pretrained model to speed up convergence. We updated batch normalization parameters during fine-tuning because of the data distribution difference between the ImageNet dataset and the Camelyon dataset. We used the Adam optimization method with a learning rate $1e-4$. The network was trained for approximately 2 epoch (500K iteration) with a batch size 32 per GPU.

To find hyperparameters and validate performance, we split Camelyon16 **train** set into our train/val set, 80% for train and 20% for validation. For AUC evaluation, we used maximum confidence probability in WSI. For FROC evaluation, we followed connected component approach [11] which find connected components and then report maximum confidence probability’s location within the component. After hyperparameter tuning, we finally train CNN with all given training dataset in Table 4.

Table 5. Feature components for predicting lymph node metastasis type.

No.	Feature description	No.	Feature description
1	largest region’s major axis length	7	maximum confidence probability in WSI
2	largest region’s maximum confidence probability	8	average of all confidence probability in WSI
3	largest region’s average confidence probability	9	number of regions in WSI
4	largest region’s area	10	sum of all foreground area in WSI
5	average of all region’s averaged confidence probability	11	foreground and background area ratio in WSI
6	sum of all region’s area		

Lymph Node Classification Module We generated the tumor probability heatmap from WSI using the metastasis detection module. For the post-processing, we thresholded the heatmap with a threshold of $t = 0.9$. We found hyperparameters and feature designs for random forest classifier in Camelyon17 **train-L** set with 5-fold cross-validation setting. Finally, we extracted 11 features described in Table 5. We built a random forest classifier to discriminate lymph node classes using extracted features. Each patient’s pN-stage was determined by the given rule [7] with the 5 lymph node slide prediction result.

Table 4. Number of training WSIs for metastasis detection module.

Training data	# of tumor slides	# of normal slides
Camelyon16 train	110	160
Camelyon16 test	50	80
Camelyon17 train-M	50*	110

* only 50 slides include region annotations from total 105 tumor slides in Camelyon17 **train-M** set

3.4 Results

Metastasis Detection on Camelyon16 We validated our metastasis detection module on the Camelyon16 dataset. For the fair comparison with the state-of-the-art methods, our model is trained on the 270 WSIs from Camelyon16 **train** set and evaluated on the 130 WSIs from Camelyon16 **test** set using the same evaluation metrics provided by the Camelyon16 challenge. Table 6 summarizes slide-level AUC and lesion-level FROC comparisons with the best previous methods. Our metastasis detection module achieved highly competitive AUC (0.9853) and FROC (0.8552) without bells and whistles.

Table 6. Metastasis detection results on Camelyon16 **test** set

Method	Ensemble	AUC	FROC
Lunit Inc.		0.985	0.855
Y. Liu et al. ensemble-of-3 [5]	✓	0.977	0.885
Y. Liu et al. 40X [5]		0.967	0.873
Harvard & MIT [11]	✓	0.994	0.807
Pathologist* [6]	-	0.966	0.724

* expert pathologist who assessed without a time constraint

Table 7. Top-10 pN-stage classification result on the Camelyon17 leaderboard [7]. The kappa score is evaluated by the Camelyon17 organizers. Accessed: 2018-03-02.

Team	Affiliation	Kappa score
Lunit Inc.*	Lunit Inc.	0.9203
HMS-MGH-CCDS	Harvard Medical School, Mass. General Hospital, Center for Clinical Data Science	0.8958
DeepBio*	Deep Bio Inc.	0.8794
VCA-TUe	Electrical Engineering Department, Eindhoven University of Technology	0.8786
JD*	JD.com Inc. - PCL Laboratory	0.8722
MIL-GPAT	The University of Tokyo, Tokyo Medical and Dental University	0.8705
Indica Labs	Indica Labs	0.8666
chengshenghua*	Huazhong University of Science and Technology, Britton Chance Center for Biomedical Photonics	0.8638
Mechanomind*	Mechanomind	0.8597
DTU	Technical University of Denmark	0.8244

* Submitted result after reopening the challenge

pN-stage Classification on Camelyon17 For validation, we first evaluated our framework on Camelyon17 **train-L** set with 5-fold cross-validation setting. Our framework achieved 0.9351 slide-level lymph node classification accuracy and 0.9017 patient-level kappa score using single CNN model in metastasis detection module. We trained additional CNN models with different model hyperparameters and fine-tuning setting. Finally, three model was ensemble by averaging probability heatmap and reached 0.9390 slide-level accuracy and 0.9455 patient-level kappa score with the 5-fold cross-validation.

Next, we evaluated our framework on the Camelyon17 **test** set and the kappa score has reached 0.9203. As shown in Table 7, our proposed framework significantly outperformed the state-of-the-art approaches by large-margins where it achieves better performance than the previous winning method (HMS-MGH-CCDS) of the Camelyon17 challenge.

Furthermore, the accuracy of our algorithm not only exceeded that of current leading approaches (bold black color in Table 8) but also significantly reduced false-negative results (red color in Table 8). This is remarkable from a clinical perspective, as false-negative results are most critical, likely to affect patient survival due to consequent delay in diagnosis and appropriate timely treatment.

Table 8. Slide-level lymph node classification confusion matrix comparison on the Camelyon17 test set. The confusion matrix is generated by the Camelyon17 organizers.

		Predicted			
		Negative	ITC	Micro	Macro
Reference	Negative	96.15%	3.08%	0.77%	0.00%
	ITC	55.88%	11.76%	32.35%	0.00%
	Micro	9.64%	2.41%	85.54%	2.41%
	Macro	3.25%	0.00%	5.69%	91.06%

(a) Ours

		Predicted			
		Negative	ITC	Micro	Macro
Reference	Negative	95.38%	0.38%	4.23%	0.00%
	ITC	76.47%	14.71%	8.82%	0.00%
	Micro	13.25%	1.20%	78.31%	7.23%
	Macro	1.63%	0.00%	12.20%	86.18%

(b) HMS-MGH-CCDS

4 Conclusion

We have introduced a robust and effective method to predict pN-stage from lymph node histological slides, using CNN based metastasis detection and random forest based lymph node classification. Our proposed method achieved the state-of-the-art result on the Camelyon17 dataset. In future work, we would like to build an end-to-end learning framework for pN-stage prediction from WSIs.

References

1. Sobin, L.H., Gospodarowicz, M.K., Wittekind, C.: TNM classification of malignant tumours. John Wiley & Sons (2011)
2. Saadatmand, S., et al.: Influence of tumour stage at breast cancer detection on survival in modern times: population based study in 173 797 patients. *Bmj* **351** (2015) h4901
3. 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. (2016) 770–778
4. Paeng, K., Hwang, S., Park, S., Kim, M.: A unified framework for tumor proliferation score prediction in breast histopathology. In: Deep Learning in Medical Image Analysis and Multimodal Learning for Clinical Decision Support. Springer (2017) 231–239
5. Liu, Y., et al.: Detecting cancer metastases on gigapixel pathology images. arXiv preprint arXiv:1703.02442 (2017)
6. Bejnordi, B.E., et al.: Diagnostic assessment of deep learning algorithms for detection of lymph node metastases in women with breast cancer. *Jama* **318**(22) (2017) 2199–2210
7. : Camelyon 2017. <https://camelyon17.grand-challenge.org/> Accessed: 2018-03-02.
8. Otsu, N.: A threshold selection method from gray-level histograms. *IEEE transactions on systems, man, and cybernetics* **9**(1) (1979) 62–66
9. Ronneberger, O., Fischer, P., Brox, T.: U-net: Convolutional networks for biomedical image segmentation. In: International Conference on Medical Image Computing and Computer-Assisted Intervention, Springer (2015) 234–241
10. Breiman, L.: Random forests. *Machine learning* **45**(1) (2001) 5–32
11. Wang, D., Khosla, A., Gargeya, R., Irshad, H., Beck, A.H.: Deep learning for identifying metastatic breast cancer. arXiv preprint arXiv:1606.05718 (2016)