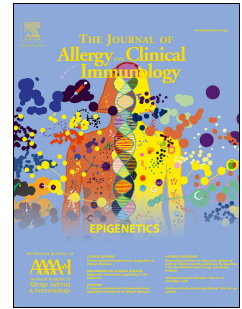# Journal Pre-proof

Expert-level Diagnosis of Nasal Polyps Using Deep Learning on Whole-slide Imaging

Qingwu Wu, MD, Jianning Chen, MD, Huiyi Deng, MD, Yong Ren, MSc, Yueqi Sun, MD, PhD, Weihao Wang, MD, Lianxiong Yuan, MD, Haiyu Hong, MD, PhD, Rui Zheng, MD, Weifeng Kong, MD, Xuekun Huang, MD, PhD, Guifang Huang, BEng, Lunji Wang, BEng, Yana Zhang, MD, PhD, Lanqing Han, MEng, Qintai Yang, MD, PhD

Please cite this article as: Wu Q, Chen J, Deng H, Ren Y, Sun Y, Wang W, Yuan L, Hong H, Zheng R, Kong W, Huang X, Huang G, Wang L, Zhang Y, Han L, Yang Q, Expert-level Diagnosis of Nasal Polyps Using Deep Learning on Whole-slide Imaging, *Journal of Allergy and Clinical Immunology* (2020), doi: https://doi.org/10.1016/j.jaci.2019.12.002.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

**Expert-level Diagnosis of Nasal Polyps Using Deep Learning on Whole-slide Imaging**

Qingwu Wu, MD[1, *], Jianning Chen, MD[2, *], Huiyi Deng, MD[1, *], Yong Ren, MSc[3, *],

Yueqi Sun, MD, PhD[4], Weihao Wang, MD[1], Lianxiong Yuan, MD[5], Haiyu Hong,

MD, PhD[6], Rui Zheng, MD[1], Weifeng Kong, MD[1], Xuekun Huang, MD, PhD[1],

Guifang Huang, BEng[3], Lunji Wang, BEng[3], Yana Zhang, MD, PhD[7], Lanqing Han,

MEng[3, #], Qintai Yang, MD, PhD[1, #]

[1] Department of Otorhinolaryngology-Head and Neck Surgery, The Third Affiliated

Hospital of Sun Yat-sen University, Guangzhou 510630, China

[2] Department of Pathology, The Third Affiliated Hospital of Sun Yat-sen University,

Guangzhou 510630, China

[3] Artificial Intelligence Innovation Center, Research Institute of Tsinghua, Pearl River

Delta, Guangzhou 510735, China

[4] Otorhinolaryngology Hospital, The First Affiliated Hospital of Sun Yat-sen

University, Guangzhou 510080, China

[5] Department of Science and Research, The Third Affiliated Hospital of Sun Yat-sen

University, Guangzhou 510630, China

[6] Department of Otolaryngology-Head and Neck Surgery, The Fifth Affiliated

Hospital of Sun Yat-sen University, Zhuhai 519020, China

[7] Feinberg School of Medicine, Northwestern University, Chicago 60611, United

21    States

22    * The authors contributed equally to this work.

23    # Corresponding authors

24    1. Lanqing Han, MEng

25    Artificial Intelligence Innovation Center

26    Research Institute of Tsinghua, Pearl River Delta

27    No. 11 Kaiyuan Road, Guangzhou 510735, China

28    Phone: 020-22213618

29    E-mail: hanlance@tsinghua-gd.org

30    2. Qintai Yang, MD, PhD

31    Department of Otolaryngology-Head and Neck Surgery

32    The Third Affiliated Hospital of Sun Yat-sen University

33    No. 600 Tianhe Road, Guangzhou 510630, China

34    Phone: 020-85252239

35    E-mail: yang.qt@163.com

36    **Funding**

39  Guangzhou (No. 201704030046), Sun Yat-sen University Clinical Research 5010

40  Program (NO.2019006) and The Third Affiliated Hospital of Sun Yat-Sen University,

41  Clinical Research Program (No. QHJH201901).

42  **Conflicts of interest**

43  The authors declare that they have no conflicts of interest. This work described was

44  original research that has not been published previously, and not under consideration

45  for publication elsewhere.

46  **Capsule summary**

47  AICEP is the first use of deep learning in combination with WSI in nasal polyp

48  diagnosis and treatment. It can improve the diagnosis and management of nasal

49  polyps more quickly and accurately.

50  **Key words**

51  CRSwNP; deep learning; pathological classification; eosinophils; WSI

52  **Acknowledgements**

53  The authors would like to thank Chunkui Shao (Professor, Department of Pathology,

54  The Third Affiliated Hospital of Sun Yat-sen University) and his colleagues for the

55  help.

56  **Ethical approval**

57  This study was approved by the Research Ethics Committee of the Institute of Basic

58  Research in Clinical Medicine, Third Affiliated Hospital of Sun Yat-sen University

59   ([2019]02-157-01). The research was registered at Chinese Clinical Trails Registry

60   (http://www.chictr.org.cn/index.aspx) with the number ChiCTR1900021601.

61    To the Editor:

62    Chronic rhinosinusitis (CRS) is defined as a chronic inflammation of the nose and

63    paranasal sinuses. It is estimated that CRS affects more than 100 million patients

64    worldwide and it involves high management costs and poor quality of life (QOL) in

65    affected subjects[1]. The presence of eosinophils in nasal polyps is linked to higher

66    postoperative visual analogue pain scores (VAS), impaired QOL, and high recurrence

67    rate[2]. A better understanding of the ratio of eosinophils (RE) to infiltrating

68    inflammatory cells in tissue is needed to improve diagnostic and treatment strategies

69    for affected patients[3].

70    Thus far, there are no uniform standards or rules regarding diagnosis of eosinophilic

71    chronic rhinosinusitis with nasal polyps (eCRSwNP), and a variety of problems exist

72    in practice. Some researchers recommend that the amounts of eosinophils per high

73    power field (HPF) should be more than 15 or 100[4, 2]. Most researchers support the

74    assessment of RE in several random HPFs, with eCRSwNP diagnosed when RE

75    is >10%[5, 6]. The traditional method ($RE_{slide-tm}$) dictates that the pathologist assesses the

76    ratio of eosinophils to infiltrating inflammatory cells (which include eosinophils,

77    neutrophils, lymphocytes, plasma cells, etc.) in 10 random HPFs for the tissue[6].

78    However, RE obviously differs between various HPFs. Preliminary studies have

79    shown sampling errors among the estimates based on 10 random HPFs and in the

80    overall eosinophil counts in the total sample. Therefore, we considered the RE of

81    whole-slide imaging (WSI) as the gold standard ($RE_{slide-actual}$) for assessing eCRSwNP

82    for its lack of sampling error. However, it is difficult in practice because it is both

83    time-consuming and subjective.

84    Artificial intelligence (AI), especially deep learning algorithms, has made great

85    progress and is similar to or even better than humans in terms of visual perception and

86    speech recognition. Therefore, we aimed to establish an artificial intelligence

87    evaluation platform (AICEP, $RE_{slide-predict}$) to diagnose eCRSwNP rapidly and

88    accurately via deep learning and WSI.

89    A total of 195 nasal polyp specimens were obtained from three affiliated hospitals of

90    Sun Yat-sen University (The Third Hospital=179, The First Hospital=9, The Fifth

91    Hospital=7). After WSI, we automatically extracted 26589 patches in the lamina

92    propria of mucosa and marked the RE in each patch ($RE_{patch-actual}$, see the Methods

93    section in this article's Online Repository at www.jacionline.org). The patches were

94    classified as the training dataset, the internal validation dataset and independent

95    external test dataset (Fig. E1).

96    In this study, our AICEP compared three common architectures (Resnet50, Xception,

97    and Inception V3) for application of a transfer learning algorithm to assess their

98    performance in the classification and regression of patches extracted from WSIs (Fig.

99    1). Within 100 epochs (iterations through the entire training dataset), the retrained

100   weights were saved due to the absence of further improvement in the mean absolute

101   error (MAE) (Fig. E2, A) and mean square error loss (Fig. E2, B).

102   First, we completed the qualitative classification of both internal validation and

6

103 external test datasets using Resnet50, Xception, and InceptionV3 models. WSI results

104 were classified as eosinophilic when $RE_{slide}$ exceeded 10% (see the Methods section

105 in this article's Online Repository at www.jacionline.org). The respective sensitivities

106 for the internal and external datasets were 97.0% and 93.5% for Resnet50, 90.1% and

107 84.2% for Xception, and 93.9% and 90.3% for InceptionV3 model, respectively. The

108 corresponding specificities were 86.0% and 84.6%, 88.2% and 88.4%, and 88.2% and

109 86.4%, individually. Our study showed that internal authentication was far superior to

110 external authentication (Fig. E3). The AUCs from internal validation and external test

111 datasets of Inception V3 were 0.974 and 0.957, respectively, which indicated that this

112 was the best model (Fig. 2, A and B).

113 Second, the convolutional neural network was visualized to identify the region of

114 eosinophils, which confirmed that the model was able to learn from the characteristics

115 of eosinophils only (Fig. 2, C and D).

116 In addition, for the quantitative analysis of AICEP, we found that the MAEs of

117 $RE_{patch-actual}$ and $RE_{patch-predict}$ in both internal validation dataset and external test

118 dataset were 4.3% and 5.8%. Meanwhile, both the consistency of intraclass

119 correlation coefficient and the agreement of $RE_{patch-predict}$ and $RE_{patch-actual}$ in the

120 internal validation dataset and external test dataset were greater than 0.95, indicating

121 high consistency from AICEP analysis (Table E1).

122 When compared with $RE_{slide-predict}$ from AICEP, pathologist simulation and $RE_{slide-actual}$

123 from the internal validation dataset of 12 patients, AICEP can diagnose all the 12

124     patients correctly, while the traditional method only made 10 correct diagnoses,

125     unfortunately, with two misdiagnosed patients (NO. 4 and 5; Fig. 2, E). Similarly,

126     when we compared $RE_{slide-predict}$ from AICEP with pathologist simulation and

127     $RE_{slide-actual}$ from the external test dataset of 16 patients, AICEP correctly diagnosed

128     all 16 patients, while the traditional method may misdiagnosed 4 patients (NO. 6, 7, 8,

129     and 10; Fig. 2, F).

130     Finally, we compared the diagnostic time between AICEP and pathologist judgement.

131     The result showed that AICEP (5.4 ± 0.87 min) took less time than $RE_{slide-tm}$ (12.7 ±

132     2.78 min) and $RE_{slide-actual}$ (148.6 ± 34.36 min, $P < 0.0001$, Table E2).

133     In our study, we advocated WSI assessment instead of $RE_{slide-tm}$. While WSI is

134     undoubtedly more accurate, it costs an immense amount of time. What's worse, in

135     China, the medical resources in the Midwest are significantly worse than those in the

136     eastern coastal areas, and pathologists are inadequate, especially in some primary

137     hospitals. To some extent, AICEP can well solve this problem, as it can diagnose

138     nasal polyp pathological types by WSI and AI more efficiently.

139     AI-facilitated diagnosis can alleviate doctors' workload and contribute to high-quality

140     medical care provision to patients in need[7, 8]. It is well known that the diagnosis of

141     disease depends on the intuition and experience of pathologists. Moreover, large

142     workload can lead to pathologists' working inefficiency and increasing the chance of

143     making mistakes. Our results showed that $RE_{slide-tm}$ may result in wrong diagnosis,

144     especially when the proportion of tissue eosinophils was approximately 10%.

145    However, this problem can be resolved by our AICEP, which can diagnose all

146    patients accurately.

147    Although AI has already shown great potentials for assisting doctors in diagnosis and

148    decision making, there are still some limitations. For instance, the real-world

149    diagnostic accuracy of AI was lower than that reported in their previous study

150    conducted with screening datasets[9]. Our study showed the similar result that AICEP

151    performed better in the internal validation dataset than in the external test dataset. In

152    our study, the internal validation dataset and training dataset came from a similar

153    process regarding slicing, staining, and WSI scanning, whereas these aspects may

154    differ in the external test dataset. Thus, it is important to optimize AICEP with data

155    from multiple centers.

156    Overall, AICEP is the first use of deep learning in combination with WSI in nasal

157    polyp diagnosis. It can evaluate the pathological characterizations of nasal polyps in a

158    faster and more accurate way. We believe that AICEP will be used widely in

159    particular in primary hospitals, even all around the world through the cloud platform.

160

161

162

163

164

Qingwu Wu, MD[1,*], Jianning Chen, MD[2,*], Huiyi Deng, MD[1,*], Yong Ren, MSc[3,*],

Yueqi Sun, MD, PhD[4], Weihao Wang, MD[1], Lianxiong Yuan, MD[5], Haiyu Hong,

MD, PhD[6], Rui Zheng, MD[1], Weifeng Kong, MD[1], Xuekun Huang, MD, PhD[1],

Guifang Huang, BEng[3], Lunji Wang, BEng[3], Yana Zhang, MD, PhD[7], Lanqing Han,

MEng[3,#], Qintai Yang, MD, PhD[1,#]

[1] Department of Otorhinolaryngology-Head and Neck Surgery, The Third Affiliated

Hospital of Sun Yat-sen University, Guangzhou 510630, China

[2] Department of Pathology, The Third Affiliated Hospital of Sun Yat-sen University,

Guangzhou 510630, China

[3] Artificial Intelligence Innovation Center, Research Institute of Tsinghua, Pearl River

Delta, Guangzhou 510735, China

[4] Otorhinolaryngology Hospital, The First Affiliated Hospital of Sun Yat-sen

University, Guangzhou 510080, China

[5] Department of Science and Research, The Third Affiliated Hospital of Sun Yat-sen

University, Guangzhou 510630, China

[6] Department of Otolaryngology-Head and Neck Surgery, The Fifth Affiliated

Hospital of Sun Yat-sen University, Zhuhai 519020, China

[7] Feinberg School of Medicine, Northwestern University, Chicago 60611, United

States

* The authors contributed equally to this work.

185    # Corresponding authors: Lanqing Han and Qintai Yang.

186    E-mail: hanlance@tsinghua-gd.org and yang.qt@163.com.

187    **References**

188    1.   Fokkens WJ, Lund VJ, Mullol J, Bachert C, Alobid I, Baroody F, et al. European

189    Position Paper on Rhinosinusitis and Nasal Polyps 2012. Rhinology Supplement 2012;

190    23: 3 p preceding table of contents, 1-298.

191    2.   Ikeda K, Shiozawa A, Ono N, Kusunoki T, Hirotsu M, Homma H, et al.

192    Subclassification of chronic rhinosinusitis with nasal polyp based on eosinophil and

193    neutrophil. The Laryngoscope 2013; 123: E1-9.

194    3.   Snidvongs K, Lam M, Sacks R, Earls P, Kalish L, Phillips PS, et al. Structured

195    histopathology profiling of chronic rhinosinusitis in routine practice. International

196    forum of allergy & rhinology 2012; 2: 376-385.

197    4.   Wen W, Liu W, Zhang L, Bai J, Fan Y, Xia W, et al. Increased neutrophilia in

198    nasal polyps reduces the response to oral corticosteroid therapy. The Journal of

199    allergy and clinical immunology 2012; 129: 1522-1528.e1525.

200    5.   Mahdavinia M, Suh LA, Carter RG, Stevens WW, Norton JE, Kato A, et al.

201    Increased noneosinophilic nasal polyps in chronic rhinosinusitis in US

202    second-generation Asians suggest genetic regulation of eosinophilia. The Journal of

203    allergy and clinical immunology 2015; 135: 576-579.

204    6.   Cao PP, Li HB, Wang BF, Wang SB, You XJ, Cui YH, et al. Distinct

205    immunopathologic characteristics of various types of chronic rhinosinusitis in adult

206    Chinese. The Journal of allergy and clinical immunology 2009; 124: 478-484,

207    484.e471-472.

208    7.   Luo H, Xu G, Li C, He L, Luo L, Wang Z, et al. Real-time artificial intelligence

209 for detection of upper gastrointestinal cancer by endoscopy: a multicentre,

210 case-control, diagnostic study. The Lancet Oncology 2019.

211 8. Lin H, Li R, Liu Z, Chen J, Yang Y, Chen H, et al. Diagnostic Efficacy and

212 Therapeutic Decision-making Capacity of an Artificial Intelligence Platform for

213 Childhood Cataracts in Eye Clinics: A Multicentre Randomized Controlled Trial.

214 EClinicalMedicine 2019; 9: 52-59.

215 9. Long E, Lin H, Liu Z, Wu X, Wang L, Jiang J, et al. An artificial intelligence

216 platform for the multihospital collaborative management of congenital cataracts.

217 Nature Biomedical Engineering 2017; 1.

218

## Figure Legend

**Figure 1.** Schematic of processes. A, nasal endoscopic examination. B, samples of nasal polyps obtained by functional endoscopic sinus surgery (FESS). C, made HE slides. D, digitized the slides into the whole slide images (WSI) by scanner. E, delineated the lamina propria to obtain region of interest (ROI). F, patches extracted from ROI of WSI and corresponding RE tags. G, trained transfer learning models that can be deployed to diagnose eCRSwNP. H, $RE_{slide-predict}$ of patients according to the model. I, chose the appropriate treatment strategy.

**Figure 2.** Performance of AICEP. A and B, the receiver operating characteristic curves (ROC) and the area under ROC (AUC) for detection of patches with RE≥10% from patches with RE<10%. A, comparison of AUC/ROC for Resnet50, Xception and Inception V3 models using internal test dataset. The Inception V3 model had an AUC (0.974) significantly greater than the other two models. B, comparison of AUC/ROC for Resnet50, Xception and Inception V3 models in independent external test dataset. The Inception V3 model also provided the best AUC (0.957) compared to the other ones. C and D, visualization and explainability of CNN models using Grad-CAM to classify patches with $RE_{patch}$≥10% from patches with $RE_{patch}$<10%. C, $RE_{patch}$=86.66%, eosinophils were marked by red arrows. D, corresponding Grad-CAM images, the highlighted areas were discriminative features of eosinophils. E and F, diagnostic efficiency comparison of AICEP and current method. Black dot represented current method result, and 50 times bootstrap were performed to evaluate its random error, blue line was the actual value of WSI and yellow line was the

14

241    AICEP predicted value of WSI, red dashed line was the diagnostic boundary of 10%.

242    E, internal validation dataset: all patients were accurately diagnosed by AICEP while

243    current method may make wrong diagnosis in patient NO. 4 and 5. F, external test

244    dataset: all patients were accurately diagnosed by AICEP while current method may

245    make wrong diagnosis in patient NO. 6, 7, 8 and 10.

**Table E1.** Consistency assessment for AICEP in internal validation dataset and

external test dataset according to the $RE_{patch-actual}$ and $RE_{slide-actual.}$

| Level | Internal Validation Dataset | | External Test Dataset | |
|---|---|---|---|---|
| | ICC Consistency | ICC Agreement | ICC Consistency | ICC Agreement |
| $RE_{patch}$ | 0.981 | 0.981 | 0.977 | 0.976 |
| | （0.979,0.983） | （0.979,0.982） | （0.975,0.979） | （0.970,0.980） |
| $RE_{slide}$ | 0.999 | 0.999 | 0.995 | 0.993 |
| | （0.997,1.000） | （0.998,1.000） | （0.985,0.998） | （0.973,0.998） |

RE, the ratio of eosinophils; ICC, intraclass correlation coefficient.

**Table E2.** Comparison of time-consuming between AICEP and pathologists.

| Method | Mean time ± SD (min) | 95% CI |
|---|---|---|
| $RE_{slide-predict}$ | 5.4±0.87 | [5.28, 5.52] |
| $RE_{slide-tm}$ | 12.7±2.78 | [12.31, 13.09] |
| $RE_{slide-actual}$ | 148.6±34.36 | [143.78, 153.42] |

| A. Identification | B. FESS | C. Slide | D. Digitization of WSI |
|---|---|---|---|



NP: nasal polyps; MT: middle turbinate; UP: uncinate process.

Figure 1,A-D

Figure 1,E-I

A

**ROC on Internal Validation Dataset**



B

**ROC on External Test Dataset**

Figure 2, A-B

C

D



Figure 2, C-D

E

**Result of Internal Validation Dataset**

F

**Result of External Test Dataset**
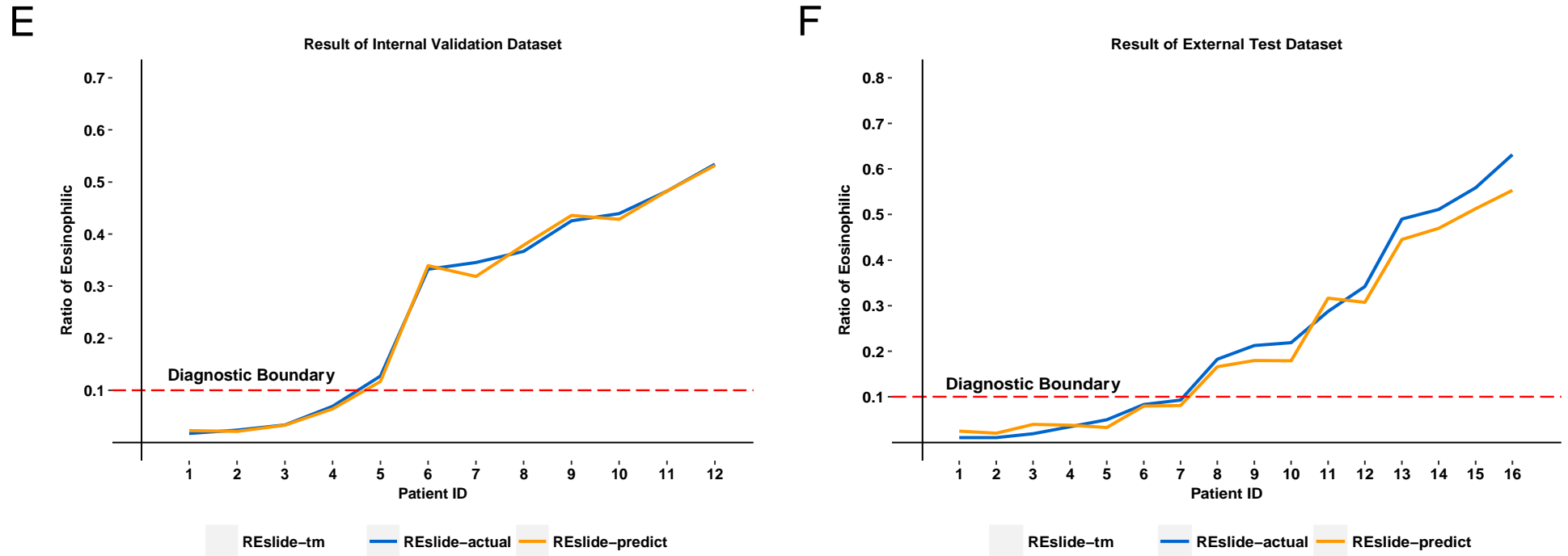
Figure 2, E-F

Figure E1

A



B



Figure E2

A

Confusion Matrix for Resnet50

| | REpatch<10% | REpatch≥10% |
|---|---|---|
| REpatch<10% | 596 | 97 |
| REpatch≥10% | 26 | 858 |

True label / Predicted label

B

Confusion Matrix for Xception

| | REpatch<10% | REpatch≥10% |
|---|---|---|
| REpatch<10% | 611 | 82 |
| REpatch≥10% | 80 | 804 |

True label / Predicted label

C

Confusion Matrix for Inception V3

| | REpatch<10% | REpatch≥10% |
|---|---|---|
| REpatch<10% | 611 | 82 |
| REpatch≥10% | 54 | 830 |

True label / Predicted label

D

Confusion Matrix for Resnet50

| | REpatch<10% | REpatch≥10% |
|---|---|---|
| REpatch<10% | 752 | 137 |
| REpatch≥10% | 70 | 1005 |

True label / Predicted label

E

Confusion Matrix for Xception

| | REpatch<10% | REpatch≥10% |
|---|---|---|
| REpatch<10% | 786 | 103 |
| REpatch≥10% | 170 | 905 |

True label / Predicted label

F

Confusion Matrix for Inception V3

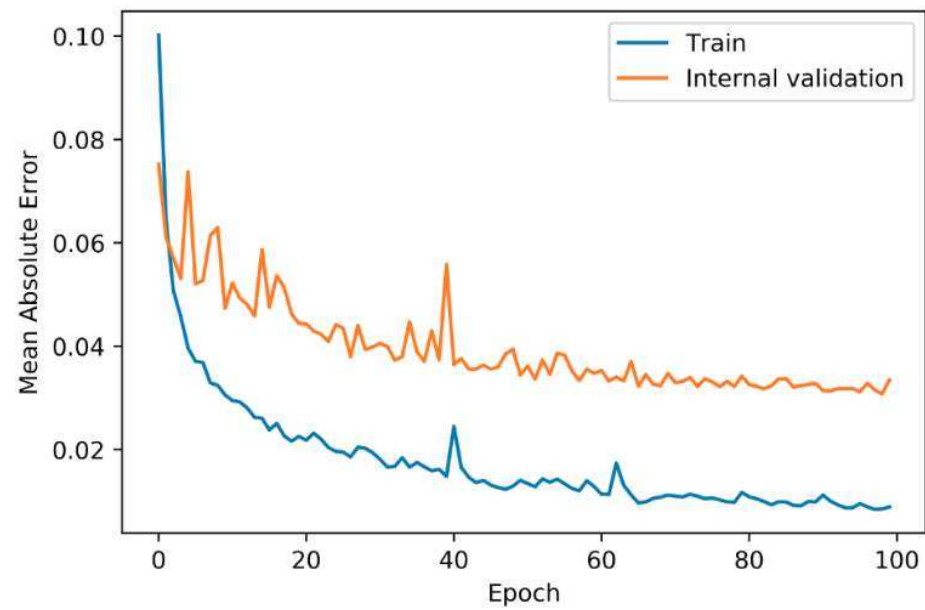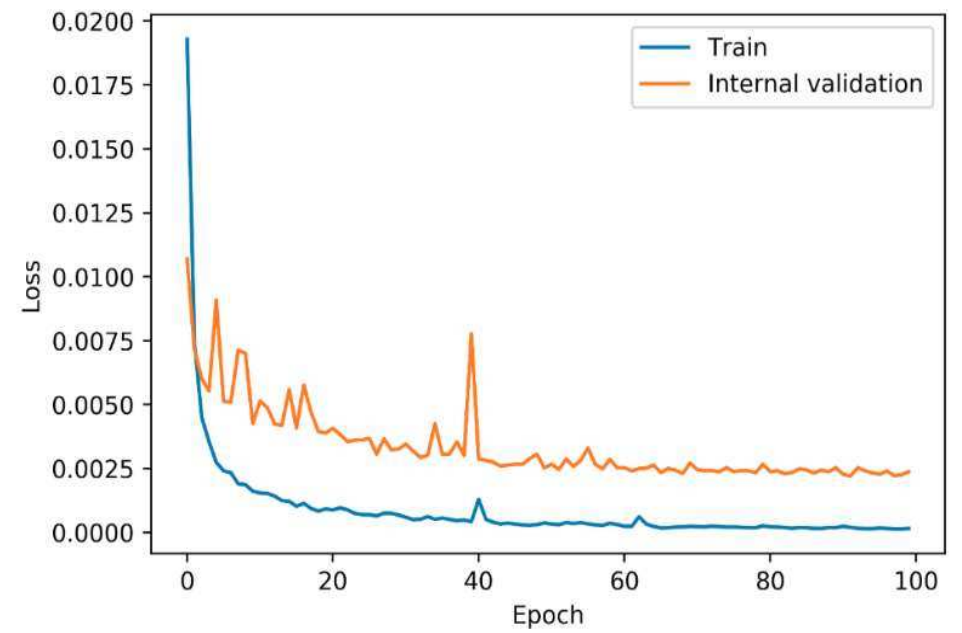| | REpatch<10% | REpatch≥10% |
|---|---|---|
| REpatch<10% | 768 | 121 |
| REpatch≥10% | 104 | 971 |

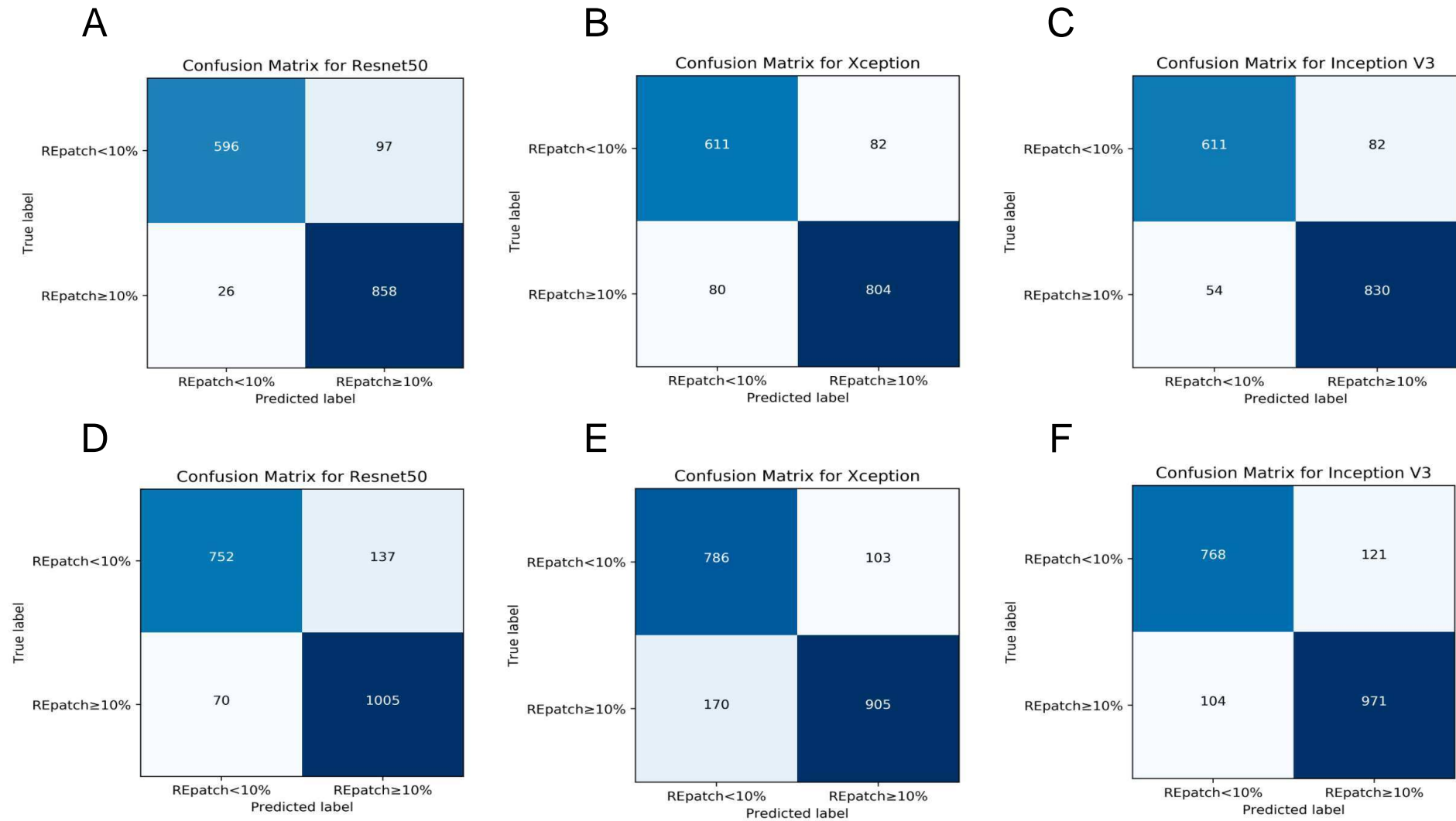True label / Predicted label

Figure E3

1  **Article's Online Repository at www.jacionline.org**

2  **METHODS**

3  **Training and internal validation datasets**

4  Biopsies of patients with CRSwNP (n = 1465) were obtained from the Department of

5  Otolaryngology in the Third Affiliated Hospital of Sun Yat-sen University (SYSU) in China

6  from January 2008 to December 2018. Following screening for staining, size, and quality of

7  specimens, 179 patients were used in this analysis. The patients were randomly divided into two

8  groups: 167 patients in the training dataset and 12 patients in the internal validation dataset. After

9  all slides were scanned through an automatic digital slide scanner (Panoramic 250 FLASH,

10  3DHISTECH Ltd., Budapest, Hungary), we obtained 179 digital whole slide images (WSIs). The

11  lamina propria of mucosa were sketched, excluding large glands, through an automated slide

12  analysis platform (ASAP) (Radboud University Medical Center, The Netherlands) to yield

13  regions of interest (ROI). Patches in ROI were automatically extracted under 400X high-power

14  field using Openslide (version 3.4.1, University of Pittsburgh, Pittsburgh, PA, USA). There were

15  167 WSIs containing 23048 patches for the training dataset and 12 WSIs containing 1577

16  patches for the internal validation dataset (Fig. E1).

17  **External test dataset**

18  Sixteen patients (16 WSIs) with nasal polyps were randomly selected from the First Affiliated

19  Hospital of SYSU (n=9) and the Fifth Affiliated Hospital of SYSU (n=7) from January 2017 to

20  December 2018. Independent preparations by each hospital were used for hematoxylin and eosin

21  staining as well as WSI scanning. In total, 1964 patches were obtained using the same method

22  mentioned above.

23    **Labeling**

24    In total, 26,589 patches were independently described and labeled by a committee comprising

25    two competent pathologists with more than 10 years of experience, and an expert pathologist

26    with more than 30 years of experience who was consulted in case of disagreement. The two

27    competent pathologists identified and counted the number of eosinophils (n1), number of

28    lymphocytes (n2), number of neutrophils (n3), and number of plasma cells (n4) in each patch.

29    The number of infiltrating inflammatory cells was regarded as the sum (t), and the ratio of

30    eosinophils ($RE_{patch-actual}$) was n1/t. When the two pathologists' assessment of $RE_{patch-actual}$

31    differed by ≤5%, the average value was used. If the difference was greater than 5%, the patch

32    was rechecked by the expert pathologist, and the value was corrected as necessary. These

33    assessments yielded the average of all patches from WSI, designated as $RE_{slide-actual}$. CRSwNP

34    patients were classified as eosinophilic when the proportion of tissue eosinophils exceeded 10%

35    of total infiltrating inflammatory cells as previously reported[1]; otherwise, they were regarded as

36    non-eosinophilic CRSwNP.

37    **Deep learning and transfer learning methods**

38    In this study, our artificial intelligence chronic rhinosinusitis evaluation platform (AICEP)

39    compared three commonly used architectures (Resnet50, Xception, and Inception V3) for

40    application of a transfer learning algorithm to assess their performance in the classification and

41    regression of patches extracted from WSIs. Each model loaded the weights pre-trained on the

42    ImageNet dataset, then removed their top layer. Next, to distinguish patches with $RE_{patch}$ values

43    greater or less than the truncated value using a classification algorithm, a full-connection layer

44    with two neurons was added and each neuron contained weights and an activation function, so it

45    can map input value to output value nonlinearly. To predict exact $RE_{patch}$ values with a regression

46  algorithm, we chose the model with the greatest area under the curve (AUC) and added a full-

47  connection layer containing only one neuron. Importantly, no activation function was used at this

48  time to ensure that the model exhibited a broader output value. Within 100 epochs (iterations

49  through the entire training dataset), the retrained weights were saved due to the absence of

50  further improvement in the mean absolute error (MAE) (Fig. E2, A) and the mean square error

51  loss (MSEL) (Fig. E2, B). Finally, the parameters of all layers of quantitative regression

52  architecture were fine-tuned in accordance with the input images and corresponding labels (Fig.

53  1).

54  To train and evaluate our models, we adopted the Keras (version 2.2) framework using

55  Tensorflow (version 1.8) backend within Python (version 3.6) programming language, including

56  libraries such as numpy, matplotlib, and Scikit-learn. Computing power was provided by one

57  Tesla V100 GPU with 32GB memory on a Nvidia DGX1 server, which had eight Tesla V100

58  GPUs, 512 GB DDR4 memory, and 7 TB SSD.

59  **Model and algorithm performance evaluation**

60  **Qualitative classification**

61  For the internal validation dataset and external test dataset using Resnet50, Xception, and

62  InceptionV3 for data training, AICEP provided an effective approach for qualitative

63  classification. WSI results were classified as eosinophilic when $RE_{slide}$ exceeded 10%, as

64  previously mentioned. The sensitivity (true positive rate) and specificity (false positive rate) of

65  the confusion matrices of these three models were calculated, as were the areas under the

66  receiver operating characteristic curve (AUC). The model with the highest AUC value was

67  selected for subsequent quantitative analyses. In addition, to verify whether the model was

68    trained correctly based on the characteristics of eosinophils, we used visual gradient-weighted

69    class activation mapping (Grad-CAM).

70    **Quantitative analysis**

71    **Evaluation of $RE_{patch}$ in internal validation and external test datasets of AICEP**

72    All patches in both internal validation and external test datasets were input into the AICEP

73    model for simulation, which produced $RE_{patch-predict}$. In addition, the MAE of $RE_{patch-predict}$ and

74    $RE_{patch-actual}$ was calculated. The concordance between $RE_{patch-predict}$ and $RE_{patch-actual}$ was evaluated

75    using the intraclass correlation coefficient.

76    **$RE_{slide}$ comparison between internal validation and external test datasets**

77    For the internal validation and external test datasets, we compared $RE_{slide-predict}$ and $RE_{slide-actual}$

78    separately. The concordance between $RE_{slide-predict}$ and $RE_{slide-actual}$ was evaluated via intraclass

79    correlation coefficient. In addition, we randomly selected 10 $RE_{patch}$ values of each WSI analysis

80    by a bootstrap method and calculated the average. The bootstrap process was repeated 50 times

81    for each WSI analysis to evaluate and compare the diagnostic effect of the traditional method

82    and of AICEP.

83    **Diagnostic time comparison between AICEP and pathologists**

84    Times for $RE_{slide-predict}$, $RE_{slide-tm}$, and $RE_{slide-actual}$ were calculated.

85    **Statistical analysis**

86    Using a bootstrap simulation of 10 random fields for diagnosis, each WSI was repeated 50 times

87    and compared with $RE_{slide-actual}$. The intraclass correlation coefficient was used to assess

88    agreement between $RE_{predict}$ with $RE_{actual}$. Receiver operating characteristic curves (ROC) were

89  adopted to evaluate the diagnostic results of AICEP on eCRSwNP. All tests were two-sided, and

90  *P*        <        0.05        was        considered        statistically        significant.

91  **References**

92  E1.    Cao PP, Li HB, Wang BF, Wang SB, You XJ, Cui YH, et al. Distinct immunopathologic

93  characteristics of various types of chronic rhinosinusitis in adult Chinese. The Journal of allergy

94  and clinical immunology 2009; 124: 478-484, 484.e471-472.

95

96
97
98

99    **Figure E1.** Workflow diagram. It illustrated the overall experimental design, and the flow of

100   whole slide images through extraction and labeling process, the training of transfer learning

101   models using internal dataset, and the evaluating of the models with internal validation dataset

102   and independent external test dataset.

103   **Figure E2.** Plot showed the performance in the training and internal validation datasets. Mean

104   absolute error was plotted against the training epoch (A) and mean square error loss was plotted

105   against the training epoch (B) during training the quantitative regression architecture over the

106   course of 100 epochs. The mean absolute error and loss of validation showed great performance

107   with little overfitting due to the diversity of the training dataset.

108   **Figure E3.** Confusion matrix of models' classification of patch with RE≥10% from patch with

109   RE<10%. A, B, C, Confusion matrix of internal validation dataset for models of Resnet50,

110   Xception and Inception V3, respectively. D, E, F, Confusion matrix of independent external test

111   dataset for models of Resnet50, Xception and Inception V3, respectively.

112

113    **Table E1.** Consistency assessment for AICEP in internal validation dataset and external test

114    dataset according to the $RE_{patch-actual}$ and $RE_{slide-actual.}$

115    **Table E2.** Comparison of time-consuming between AICEP and pathologists.

116