1 **Title: Assessment of an AI virtual staining model performance across same and serial tissue**

2 **sections using CD3+ T cell ground truth**

3

4 **Abu Bakr Azam[1†], Felicia Wee[2†], Juha P. Väyrynen[3], Willa Wen-You Yim[2], Yue Zhen**

5 **Xue[2], Bok Leong Chua[1], Jeffrey Chun Tatt Lim[2], Daniel Shao Weng Tan[4], Angela**

6 **Takano[5], Chun Yuen Chow[5], Li Yan Khor[5], Tony Kiat Hon Lim[5], Joe Yeong[2,5*], Mai Chan**

7 **Lau[6,7*], Yiyu Cai[1*]**

8 [1]School of Mechanical and Aerospace Engineering, Nanyang Technological University,

9 Singapore 639798

10 [2]Institute of Molecular and Cell Biology, Agency for Science, Technology and Research

11 (A*STAR), Singapore 138673

12 [3]Translational Medicine Research Unit, Medical Research Center Oulu, Oulu University

13 Hospital, and University of Oulu, POB 5000, 90014 Oulu, Finland

14 [4]Division of Medical Oncology, National Cancer Centre, Singapore 168583

15 [5]Department of Anatomical Pathology, Division of Pathology, Singapore General Hospital,

16 Singapore 169856

17 [6]Bioinformatics Institute (BII), Agency for Science, Technology and Research (A*STAR), 30

18 Biopolis Street, Matrix, Singapore 138671, Republic of Singapore

19 [7]Singapore Immunology Network (SIgN), Agency for Science, Technology and Research

20 (A*STAR), 8A Biomedical Grove, Immunos, Singapore 138648, Republic of Singapore.

21 †These authors contributed equally to this work and share first authorship

22

23 **\*Corresponding Authors:**

24 Joe YEONG

25    Institute of Molecular and Cell Biology (IMCB), A*STAR, Proteos, 61 Biopolis Drive,

26    Singapore 138673

27    Department of Anatomical Pathology, Division of Pathology, Singapore General Hospital,

28    Academia, 20 College Road, Singapore 169856

29    Email: yeongps@imcb.a-star.edu.sg, Tel: +65 65869527

30

31    Yiyu CAI

32    School of Mechanical and Aerospace Engineering, Nanyang Technological University, 50

33    Nanyang Avenue, Singapore 639798

34    Email: MYYCai@ntu.edu.sg, Tel: +65 67905777

35

36    Mai Chan LAU

37    Bioinformatics Institute, Agency for Science, Technology and Research (A*STAR), Matrix, 30

38    Biopolis Street, Singapore 138671

39    Singapore Immunology Network, Agency for Science, Technology and Research (A*STAR), 8a

40    Biomedical Grove, Singapore 138648

41    Email: Lau_Mai_Chan@bii.a-star.edu.sg, Tel: +65 64070606

45

46

47

48    **Abstract: (290 words/300)**

49        Immunophenotyping via multi-marker assays significantly contributes to patient selection,

50    therapeutic monitoring, biomarker discovery, and personalized treatments. Despite its potential,

51    the multiplex immunofluorescence (mIF) technique faces adoption challenges due to technical and

52    financial constraints. Alternatively, hematoxylin and eosin (H&E)-based prediction models of cell

53    phenotypes can provide crucial insights into tumor-immune cell interactions and advance

54    immunotherapy. Current methods mostly rely on manually annotated cell label ground truths, with

55    limitations including high variability and substantial labor costs. To mitigate these issues,

56    researchers are increasingly turning to digitized cell-level data for accurate in-situ cell type

57    prediction. Typically, immunohistochemical (IHC) staining is applied to a tissue section serial to

58    one stained with H&E. However, this method may introduce distortions and tissue section shifts,

59    challenging the assumption of consistent cellular locations. Conversely, mIF overcomes these

60    limitations by allowing for mIF and H&E staining on the same tissue section. Importantly, the

61    multiplexing capability of mIF allows for a thorough analysis of the tumor microenvironment by

62    quantifying multiple cell markers within the same tissue section. In this study, we introduce a

63    Pix2Pix generative adversarial network (P2P-GAN)-based virtual staining model, using CD3$^+$ T-

64    cells in lung cancer as a proof-of-concept. Using an independent CD3 IHC-stained lung cohort,

65    we demonstrate that the model trained with cell label ground-truth from the same tissue section as

66    H&E staining performed significantly better in both CD3$^+$ and CD3$^-$ T-cell prediction. Moreover,

67    the model also displayed prognostic significance on a public lung cohort, demonstrating its

68    potential clinical utility. Notably, our proposed P2P-GAN virtual staining model facilitates image-

69    to-image translation, enabling further spatial analysis of the predicted immune cells, deepening

70    our understanding of tumor-immune interactions, and propelling advancements in personalized

3

71    immunotherapy. This concept holds potential for the prediction of other cell phenotypes, including

72    CD4[+], CD8[+], and CD20[+] cells.

73

74    **Introduction**

75    Immune phenotyping in tissue, facilitated by multi-marker assays such as mIF, plays a pivotal

76    role in patient selection, treatment monitoring, biomarker discovery, and the development of

77    targeted and personalized therapeutic strategies[1,2,3]. Nevertheless, the wider adoption of the mIF

78    technique faces challenges as it remains inaccessible to many laboratories due to technical and

79    time constraints or funding limitations. Conversely, the utilization of hematoxylin and eosin

80    (H&E)-based prediction models present a viable alternative for generating data to enhance our

81    comprehension of the intricate interactions within the immune system. Given that H&E staining

82    is cost-effective and routinely performed in numerous histology laboratories, integrating H&E-

83    based prediction models into existing workflows can be achieved with relative ease. This approach

84    has the potential to revolutionize the field of immunotherapy, opening new avenues for

85    advancements in treatment strategies.

86    Current studies of H&E-based approaches largely rely on manual annotated cell label ground

87    truth[4,5]. For instance, a study by Wilde *et al.* demonstrated the use of deep learning (DL) to assess

88    two prognostic risk parameters, OP-TIL and the multinucleation index (MuNI), in hematoxylin

89    and eosin (H&E) stained slides from patients with oropharyngeal squamous cell carcinoma[6]. The

90    group proposed two DL-based imaging biomarkers, namely OP-TIL, which quantitatively

91    characterizes the spatial patterns between tumor infiltrating lymphocytes (TILs) and their

92    surrounding cells[7], while the MuNI quantifies the multinucleated tumor cells in epithelial regions[8].

93    Conditional generative adversarial network (cGAN) models were adopted for cell segmentation

94    based on OP-TIL, and trained for in-silico computation of MuNI. This group also highlighted the

95   potential clinical importance of identification and tissue localization of TIL subtypes, such as those

96   expressing CD4, CD8, and CD20[7]. However, the applicability of these approaches was limited by

97   the availability of manual annotation of TILs and multinucleated tumor cells by pathologists, with

98   high inter- and intra-observer variability and high labor costs[9].

99   To address both inter- and intra-observer discrepancies in the annotation and scoring of cell

100  phenotypes, there has been a growing interest in the utilization of digitized cell-level data as the

101  definitive reference for predicting cell types in situ[10]. Commonly, immunohistochemical (IHC)

102  staining is applied to a tissue section that is consecutive to another one stained with H&E, assuming

103  that similar cells maintain identical locations across both sections. Yet, in conventional IHC

104  methods, manual preparation can cause distortions, and heat fixation can shift the tissue section[11],

105  disrupting this assumption. Furthermore, achieving same-section ground truth is impeded by

106  chromogenic IHC due to deposition of the brown chromogen 3,3'-diaminobenzidine (DAB).

107  Alternatively, multiplex immunofluorescence (mIF) overcomes these limitations, enabling

108  staining on the same tissue section used for H&E staining. Crucially, mIF's multiplexing feature

109  allows a comprehensive analysis of the tumor microenvironment (TME) by quantifying multiple

110  cell markers within the same tissue section[12]. In the realm of immunotherapy, the simultaneous

111  quantification of immune markers like CD3, CD4, CD8, cytokeratin, PD-1, and CTLA-4 within

112  the same tissue space is critical for a comprehensive understanding of tumor-immune

113  interactions[13,14]. Here, we propose a Pix2Pix generative adversarial network (P2P-GAN)-based

114  virtual staining model, using CD3+ T-cells in lung cancer as the study model (Figure 1a). The

115  choice of CD3+ T-cells highlights their significant role in lung cancer prognosis and

116  treatment[15,16,17]. We hypothesize that the performance of the prediction model can be impacted by

117  cellular differences in adjacent (non-identical) tissue sections. To test this hypothesis, we built and

118     compared two DL models, one trained using the CD3$^+$ T-cell ground truth obtained by mIF staining

119     of the same tissue section stained with H&E (abbreviated as same-section model; Figure 1b), while

120     the other model was trained using the CD3$^+$ T-cell ground truth obtained by mIF staining of the

121     serial tissue section stained with H&E (abbreviated as serial-section model; Figure 1b).

122

123     **Materials/Subjects and Methods**

124     <u>Cohorts</u>

125     This study was conducted using three lung cancer cohorts (2 in-house and 1 public). The

126     training cohort consisted of formalin-fixed paraffin embedded (FFPE) tissues in the tissue

127     microarray (TMA) format prepared in the Department of Anatomical Pathology of Singapore

128     General Hospital (Agency of Science, Technology and Research (A*STAR) IRB: 2021-161, 2021-

129     188, 2021-112). The tissue sections were stained with H&E and mIF (anti-CD3 and DAPI for

130     nuclear staining) in the Institute of Molecular and Cell Biology (IMCB) at the Agency for Science,

131     Technology and Research, Singapore. Using this cohort, we prepared the same-section and serial-

132     section datasets. In the same-section dataset, 57 H&E and mIF image pairs were generated from

133     the same tissue sections of the 57 patients. In the serial-section dataset, a separate set of H&E

134     images were generated using tissue sections adjacent to the tissue sections used for the mIF

135     staining.

136     Separate in-house and public cohorts were used for evaluation of the model performance

137     (Table 1). The in-house cohort comprised CD3 IHC-stained images along with H&E images

138     generated from the corresponding serial-section in TMAs (designated the IHC cohort). The public

139     cohort consisted of H&E-stained images (20× magnification) and the companion patient survival

6

140    data were downloaded from OncoSG (Singapore Oncology, Data Portal) (designated the Onco-

141    SG cohort).

142    <u>Tissue staining</u>

143        The FFPE tissues were sectioned (4 µm thickness) and heat-fixed at 65°C for 5 min before

144    manual staining with hematoxylin (Epredia, Fisher Scientific, Porto Salvo, Portugal) and eosin

145    (Epredia, Fisher Scientific, Gothenburg, Sweden). IHC staining was performed on the FFPE

146    tissues (4 µm thickness) with anti-CD3 primary antibody (1:200; Dako A0452, Santa Clara, CA,

147    USA) using the Leica Bond Max autostainer (Leica Biosystems, Melbourne, Australia) and Bond

148    Refine Detection Kit (Leica Biosystems) as previously described[19]. The H&E and IHC stained

149    slides were then scanned using the Axioscan.Z1 Slide Scanner (Zeiss, Oberkochen, Germany).

150        Next, mIF staining was performed on the FFPE tissue sections (4 µm thickness) using the

151    Leica Bond Max autostainer (Leica Biosystems, Melbourne, Australia), Bond Refine Detection

152    Kit (Leica Biosystems) and Opal 6-Plex Detection Kit for Whole Slide Imaging (Akoya

153    Biosciences, Marlborough, MA, USA) as previously described[19]. In brief, FFPE tissue sections

154    were subjected to repeated cycles of heat-induced epitope retrieval, incubation with anti-CD3

155    primary antibody (Dako #A0452), anti-rabbit poly-HRP-IgG (Ready-to-use; Leica Biosystems)

156    and Opal tyramide signal amplification (TSA) (Akoya Biosciences). Spectral DAPI (4',6-

157    diamidino-2-phenylindole) (Akoya Biosciences) was applied as the final nuclear counterstain.

158    Images were captured using the Vectra 3 Automated Quantitative Pathology Imaging System

159    (Akoya Biosciences). After scanning, the mIF slides were subjected to H&E staining, followed by

160    scanning on the Axioscan.Z1 Slide Scanner (Zeiss).

161    Ground truth cell labels

162    For model training, ground truth cell labelling involved the identification of CD3+ cells in the

163    H&E image space according to a series of steps. First, nuclei in the H&E image were identified

164    using the StarDist Python library (pre-trained for H&E images)[20]. Second, nuclei and CD3+ regions

165    in the mIF image were identified individually using the StarDist Python library (pre-trained for

166    fluorescence images) based on DAPI staining and CD3 expression, respectively. These regions

167    were then overlaid to identify CD3+ T-cells in the mIF image. Third, the CD3+ T-cells identified

168    in the mIF image were matched to the closest nuclei in the H&E image stained on the same (post-

169    mIF H&E staining) or serial tissue section (designated same-section and serial-section datasets,

170    respectively). The H&E image with CD3+ information (i.e., ground truth image) was then

171    deconvoluted into red (R), green (G), and blue (B) channels representing the CD3+ T-cell,

172    haematoxylin (H), and eosin (E) staining, respectively. Representation of the CD3+ T-cell

173    information in a separate channel i.e., R, facilitates the identification of predicted CD3+ T-cells

174    during model deployment. Considering that CD3 localizes to the cell membrane whereas DAPI

175    staining is localized in the nucleus, Gaussian noise (kernel size 101) was applied to the R channel

176    of the image to increase the spread of the CD3+ signals while keeping the maximum intensity at

177    its center. This facilitates the identification of predicted CD3+ T-cells, which relies on an overlap

178    between CD3 and DAPI intensities i.e., R and G channels.

179    In the IHC testing dataset, CD3 signal localization in an IHC image was first determined by

180    applying a threshold (value >100) to the DAB stain intensity, resulting in a binary mask where 1

181    indicates CD3 detection and 0 indicates otherwise. The CD3 mask was then overlaid on the nuclei

182    segmented in the paired H&E image to identify CD3+ T-cells (ground truth cell labels) according

183    to the same procedure described for mIF dataset. In the Onco-SG testing dataset, two pathologists

184 (YZX and JPV) assessed the H&E images and scored the %TIL. Model performance was evaluated

185 by comparing overall %CD3$^+$ T-cell with the %TIL in individual patients by Spearman's

186 correlation analysis. We also assessed the 5-year overall survival association with the patient

187 groups stratified using the mean DL-predicted %CD3$^+$ T-cell versus the mean of the %TIL values

188 determined by the two pathologists. If multiple images were available for the same patient, the

189 patient-average %CD3$^+$ T-cell or %TIL value was used.

190 P2P-GAN model architecture

191  A conventional GAN incorporates a generative network to produce image candidates and a

192 discriminative network for their evaluation. The former network is trained to 'fool' the latter, hence

193 facilitating unsupervised learning by the model. The P2P-GAN is a variation of a conditional

194 GAN, in which the generator output image is conditional on the input image, and hence is designed

195 perfectly for the image-to-image translation task. In this study, we adopted the P2P-GAN

196 architecture reported by Isola *et al.*[21] in which a U-Net was used as the generator and a

197 convolutional neural network (CNN) was used as the discriminator (Figure 2). Model training

198 involved presenting the generator with stain-deconvoluted H&E images, while presenting the

199 discriminator with ground truth images (i.e., stain-deconvoluted H&E images overlaid with mIF-

200 identified CD3$^+$ T-cell information). These images were then compared with the generator

201 predicted images to output a 30×30 matrix for updating both the generator and discriminator

202 (Figure 2; more details are provided below).

203 Model training

204  Two P2P-GAN models were trained using the same-section and serial-section training

205 datasets (henceforth referred to as the same-section and serial-section models, respectively). Each

206 image in the training dataset (Table 1) was divided into 256×256 image patches (total 9,633

207   patches). Of these, 96% (9,249 patches) were used for model training and 4% (384 patches) were

208   randomly selected for model testing (hereafter referred as the held-out subset). The generator and

209   discriminator work in an adversarial fashion such that the respective losses are balanced out. The

210   overall objective is to reach an optimum for the two conflicting goals, where the generator

211   produces an output that is almost indistinguishable from the ground truth images, while the

212   discriminator can distinguish images generated by the generator from ground truth images.

213   Overall, three different types of losses must be minimized: LOSS 1, which measures the mean

214   absolute difference between the generator output image and the ground truth image, is used to

215   update the generator network; LOSS2/LOSS 3 and LOSS 4 measure the difference between the

216   $30 \times 30$ feature matrix output from the discriminator with two $30 \times 30$ target matrices, one of which

217   contains all 0 digits and the other contains all 1 digits. This allows quantification of 'lack of

218   capability' and 'capability', respectively, of the discriminator in distinguishing the generator

219   output image; LOSS2 (essentially LOSS3) is feedback to the generator, while LOSS 3 and LOSS4

220   are feedback to the discriminator (Figure 2). The training of both models involved 150 epochs with

221   a batch size of 350. A regularization value of 100 was applied to LOSS 1 (i.e., the mean absolute

222   loss).

223   Model performance characteristics

224      Model performance was quantified based on two key metrics, namely $CD3^+$ and $CD3^-$ T-cell

225   counts, and overall accuracy (defined as the ratio of correctly predicted $CD3^+$ and $CD3^-$ T-cell

226   counts to the total number of cells). The model-predicted $CD3^+$ and $CD3^-$ T-cell counts were

227   identified as shown in Figure 3. Specifically, model-predicted CD3 signals (represented in the red

228   channel) were overlaid with the nuclei segmented from the input H&E image to identify the $CD3^+$

229   T-cell, whereas nuclei (or cells) with no matching CD3 signals were deemed to be $CD3^-$ T cells.

10

230    The model-predicted CD3$^+$ and CD3$^-$ T-cell values were then overlaid with the paired mIF

231    (training cohort) or IHC (testing cohort) images to quantify the accurately predicted CD3$^+$ and

232    CD3$^-$ T-cell counts.

233

234    **Results**

235    <u>Validating model performance using training samples</u>

236        As a sanity check, we assessed the model performance with the image patches used for training

237    (Table 1; N = 57). Of note, the same-section and serial-section datasets were used for testing the

238    same-section and serial-section models, respectively. The predicted CD3$^+$ and CD3$^-$ T-cell counts

239    from both same-section and serial-section models were highly comparable to the mIF-quantified

240    CD3$^+$ and CD3$^-$ T-cell counts (i.e., ground truth; all $p < 0.005$) with Pearson's correlation $>0.95$

241    (Figure 4a-d). However, based on the Mann-Whitney U-test, the same-section model outperformed

242    the serial-section model by a slight margin in terms of overall accuracy (Figure 4e; $p < 0.005$).

243    <u>Performance comparison of same-section and serial-section models with held-out training cohort</u>

244        We randomly selected 4% of image patches (384 patches) from the same-section training

245    cohort for model testing. While model-predicted CD3$^+$ and CD3$^-$ T-cell counts from both the same-

246    section and serial-section models were reasonably comparable to the mIF-quantified CD3$^+$ and

247    CD3$^-$ T-cell counts (i.e., ground truth) (all $p < 0.005$, Figure 5), same-section model predictions

248    showed better concordance with the ground truth as compared with that of serial-section model

249    (Pearson's correlation coefficients 0.784 vs. 0.733, and 0.675 vs. 0.57, respectively; Figure 5a-d).

250    Based on Mann-Whitney U-tests, there was no significant difference in the overall accuracy of the

251    same-section and serial-section models (Figure 5e; $p = 0.62$).

11

252 Performance comparison of same-section and serial-section models on an independent IHC cohort
253 (N = 48)

254     In agreement with the results from the held-out cohort analysis, the $CD3^+$ and $CD3^-$ T-cell

255 counts predicted by both the same-section and serial-section models (Figure 6a) corresponded

256 closely to the IHC-quantified counts, representing the ground truth (p < 0.005, Figures 6b-e).

257 Importantly, the same-section model outperformed the serial-section model, displaying stronger

258 correlations with the IHC ground truth (Figure 6b-c; $CD3^-$ T-cell $r = 0.85$ vs. 0.678; Figure 6d-e;

259 $CD3^+$ T-cell $r = 0.886$ vs. 0.798), and achieving a higher average accuracy (Figure 6f; mean

260 accuracy = 0.92 vs. 0.65).

261

262 Validating the prognostic association of model-predicted $CD3^+$ T-cells

263     Evaluation of the models' performance on the public Onco-SG cohort (Figure 7a), composed

264 of 204 lung samples (Table 1), revealed a significant correlation between model-predicted CD3

265 patient groups and 5-year overall survival (Figure 7b; p = 0.013). This association was more

266 pronounced than that observed when patient stratification was based on manual TIL scoring by

267 two pathologists (Figure 7c-d; p = 0.3 and p = 0.06), suggesting the added value of our model in

268 predicting patient outcomes. Nonetheless, the abundance of model-predicted $CD3^+$ T-cells showed

269 significant correspondence with the TIL scoring by both pathologists (Figure 7e; p < 0.05).

270

271 **Discussion**

272     In this study, we developed and examined P2P-GAN virtual staining models to predict $CD3^+$

273 T-cells from low-cost digitized H&E images. A significant aspect of our investigation was the

274 exploration of performance disparities that arise when ground truth cell labels are obtained from

275 the same tissue section used for H&E staining, as opposed to a serial section. Our findings

276 demonstrate that the model trained using the same-section approach consistently surpasses the

12

277     serial-section model. This superiority manifests as stronger correlations with mIF and IHC-

278     quantified $CD3^+$ and $CD3^-$ T-cells, along with heightened overall prediction accuracies. It also

279     reinforces the potential of the same-section model as a robust technique in histopathology-driven

280     immune phenotyping. Crucially, our work also showcased the enhanced prognostic utility of our

281     model-predicted $CD3^+$ T-cell abundance when compared to traditional manual TIL scores. This

282     emphasizes the clinical relevance of our proposed virtual staining model in a real-world setting,

283     potentially facilitating improved patient stratification and treatment decision-making. A distinctive

284     feature that sets our proposed model apart from traditional DL models for cell prediction is its

285     capability for image-to-image translation, virtually staining the CD3 marker within the original

286     H&E context. This has two major implications. First, it facilitates further downstream analysis of

287     the TME and spatial interplay between predicted cell types and other cellular or tissue data derived

288     from H&E through either pathological assessment or digital pathology. Second, it creates a new

289     pathway for integrating incremental cell type predictions from different models onto the same

290     H&E space. Collectively, these advancements could significantly enhance our understanding of

291     the TME, potentially leading to the identification of novel spatial biomarkers or therapeutic targets.

292      While our proposed approach has yielded encouraging results, it is important to acknowledge

293     its inherent limitations. First, our current model is designed specifically for $CD3^+$ T-cells

294     prediction from H&E images and may not generalize well to other cell types or markers without

295     significant adjustments or retraining. Additionally, its performance may be compromised when

296     applied to tumor types beyond lung cancer. Second, the application of this model is largely limited

297     to high-quality digital slides. Its performance may be affected by variations in tissue preparation,

298     staining procedures, and image acquisition methods across different laboratories. Nevertheless, the

299     clinical significance of our model has been validated using a publicly available lung cohort. Lastly,

300 despite overall robust performance, we noted outliers in our model's predictions, indicating

301 potential areas for improvement. These discrepancies suggest complex, unaddressed variables

302 within biological samples that need further investigation. Future endeavors should focus on

303 understanding these outlier causes, refining modeling techniques, and incorporating larger, more

304 diverse datasets for improved generalizability and outlier management.

305     In conclusion, our thorough exploration into the necessity of employing ground truth cell

306 labels from identical tissue sections in a $CD3^+$ T-cell prediction model signifies a notable advance

307 in the domain of H&E-based virtual staining research. Our novel image-to-image translation

308 capability paves the way for in-depth TME analyses. Combined with the potential of predicting

309 refined cell types via the mIF technique, our model unveils exciting new possibilities for biomarker

310 discovery and the advancement of therapeutic strategies. While certain limitations are observed,

311 these challenges underscore the direction for future investigations, the results of which could

312 greatly enhance the prediction accuracy and clinical applicability of this innovative approach.

313

314 **Acknowledgements**

318

319 **Conflict of Interest**

320 All authors declare no conflict of interest.

321

322 **Ethics Approval and Consent to Participate**

14

323    This study was approved by the Agency of Science, Technology and Research (A*STAR) Human

324    Biomedical Research Office (A*STAR IRB: 2021-161, 2021-188, 2021-112).

325

326    **Author Contributions**

327    J.P.S.Y, M.C.L. and Y.C. conceived and directed the study. A.B.A. performed the development,

328    training and testing of the DL models and conducted the biostatistical analysis; B.L.C. performed

329    the testing of codes. F.W. and J.C.T.L. performed immunohistochemical techniques; J.P.V. and

330    Y.Z.X. performed the TIL scoring. W.W.Y. created the publication figures. D.S.W.T, A.T.,

331    C.C.Y., L.Y.K. and T.K.H.L. conducted the sample acquisition and provided clinical pathological

332    and oncological perspectives. A.B.A., F.W. and M.C.L. prepared the manuscript. All authors

333    reviewed the manuscript.

334

346

**Data Availability Statement**

The mIF and in-house IHC data sets used during the current study are available from the corresponding author upon reasonable request. The external lung cancer cohort is available in the OncoSG repository, https://src.gisapps.org/OncoSG/. The scripts used in this study can be found in the following GitHub repository, https://github.com/abubakrazam/Pix2Pix_TIL_H-E.git

352

**References**

1. Parra ER, Villalobos P, Behrens C, Jiang M, Pataer A, Swisher SG, et al. Effect of neoadjuvant chemotherapy on the immune microenvironment in non–small cell lung carcinomas as determined by multiplex immunofluorescence and image analysis approaches. *J ImmunoTher Cancer*. 2018;6(1):48.

2. Pilla L, Maccalli C. Immune Profiling of Cancer Patients Treated with Immunotherapy: Advances and Challenges. *Biomedicines*. 2018;6(3):76.

3. Cascone T, Sepesi B, Lin HY, Kalhor N, Parra ER, Jiang M, et al. A Phase I/II Study of Neoadjuvant Cisplatin, Docetaxel, and Nintedanib for Resectable Non–Small Cell Lung Cancer. *Clin Cancer Res*. 2020;26(14):3525-36.

4. Shamai G, Livne A, Polónia A, Sabo E, Cretu A, Bar-Sela G, et al. Deep learning-based image analysis predicts PD-L1 status from H&E-stained histopathology images in breast cancer. *Nat Commun*. 2022;13(1):6753.

366    5. Duenweg SR, Brehler M, Bobholz SA, Lowman AK, Winiarz A, Kyereme F, et al. Comparison

367    of a machine and deep learning model for automated tumor annotation on digitized whole slide

368    prostate cancer histology. *PLoS One*. 2023;18(3):e0278084.

369    6. Wilde DC, Castro PD, Bera K, Lai S, Madabhushi A, Corredor G, et al. Oropharyngeal cancer

370    outcomes correlate with p16 status, multinucleation and immune infiltration. *Mod Pathol*.

371    2022;35(8):1045-54.

372    7. Corredor G, Toro P, Koyuncu C, Lu C, Buzzy C, Bera K, et al. An Imaging Biomarker of

373    Tumor-Infiltrating Lymphocytes to Risk-Stratify Patients With HPV-Associated Oropharyngeal

374    Cancer. *J Natl Cancer Inst*. 2022;114(4):609-17.

375    8. Koyuncu CF, Lu C, Bera K, Zhang Z, Xu J, Toro P, et al. Computerized tumor multinucleation

376    index (MuNI) is prognostic in p16+ oropharyngeal carcinoma. *J Clin Invest*. 2021;131(8).

377    9. Schalper KA, Brown J, Carvajal-Hausdorf D, McLaughlin J, Velcheti V, Syrigos KN, et al.

378    Objective measurement and clinical significance of TILs in non-small cell lung cancer. *J Natl*

379    *Cancer Inst*. 2015;107(3).

380    10. Roy M, Wang F, Teodoro G, Bhattarai S, Bhargava M, Rekha TS, et al. Deep learning based

381    registration of serial whole-slide histopathology images in different stains. *J Pathol Inform*.

382    2023;14:100311.

383    11. Grizzle WE. Special symposium: fixation and tissue processing models. *Biotech Histochem*.

384    2009;84(5):185-93.

385   12. Tan WCC, Nerurkar SN, Cai HY, Ng HHM, Wu D, Wee YTF, et al. Overview of multiplex

386   immunohistochemistry/immunofluorescence techniques in the era of cancer immunotherapy.

387   *Cancer Commun (Lond)*. 2020;40(4):135-153.

388   13. Topalian SL, Hodi FS, Brahmer JR, Gettinger SN, Smith DC, McDermott DF, et al. Safety,

389   activity, and immune correlates of anti-PD-1 antibody in cancer. *N Engl J Med*.

390   2012;366(26):2443-54.

391   14. Waldman AD, Fritz JM, Lenardo MJ. A guide to cancer immunotherapy: from T cell basic

392   science to clinical practice. *Nature Reviews Immunology*. 2020;20(11):651-68.

393   15. Al-Shibli KI, Donnem T, Al-Saad S, Persson M, Bremnes RM, Busund LT. Prognostic effect

394   of epithelial and stromal lymphocyte infiltration in non-small cell lung cancer. *Clin Cancer Res*.

395   2008;14(16):5220-7.

396   16. Chen B, Li H, Liu C, Xiang X, Wang S, Wu A, et al. Prognostic value of the common tumour-

397   infiltrating lymphocyte subtypes for patients with non-small cell lung cancer: A meta-analysis.

398   *PLoS One*. 2020;15(11):e0242173.

399   17. Geng Y, Shao Y, He W, Hu W, Xu Y, Chen J, et al. Prognostic Role of Tumor-Infiltrating

400   Lymphocytes in Lung Cancer: a Meta-Analysis. *Cell Physiol Biochem*. 2015;37(4):1560-71.

401   18. cBioPortal for Cancer Genomics [cited 2023 May 21].

402   19. Lim JCT, Yeong JPS, Lim CJ, Ong CCH, Wong SC, Chew VSP, et al. An automated staining

403   protocol for seven-colour immunofluorescence of human tissue sections for diagnostic and

404   prognostic use. *Pathology*. 2018 Apr;50(3):333-341

18

405    20. Schmidt U, Weigert M, Broaddus C, Myers G, editors. Cell Detection with Star-Convex

406    Polygons. Medical Image Computing and Computer Assisted Intervention – MICCAI 2018; 2018

407    2018//; Cham: Springer International Publishing.

408    21. Isola P, Zhu J-Y, Zhou T and Efros AA. Image-to-Image Translation with Conditional

409    Adversarial Networks. Computer Science: Computer Vision and Pattern Recognition. 2018.

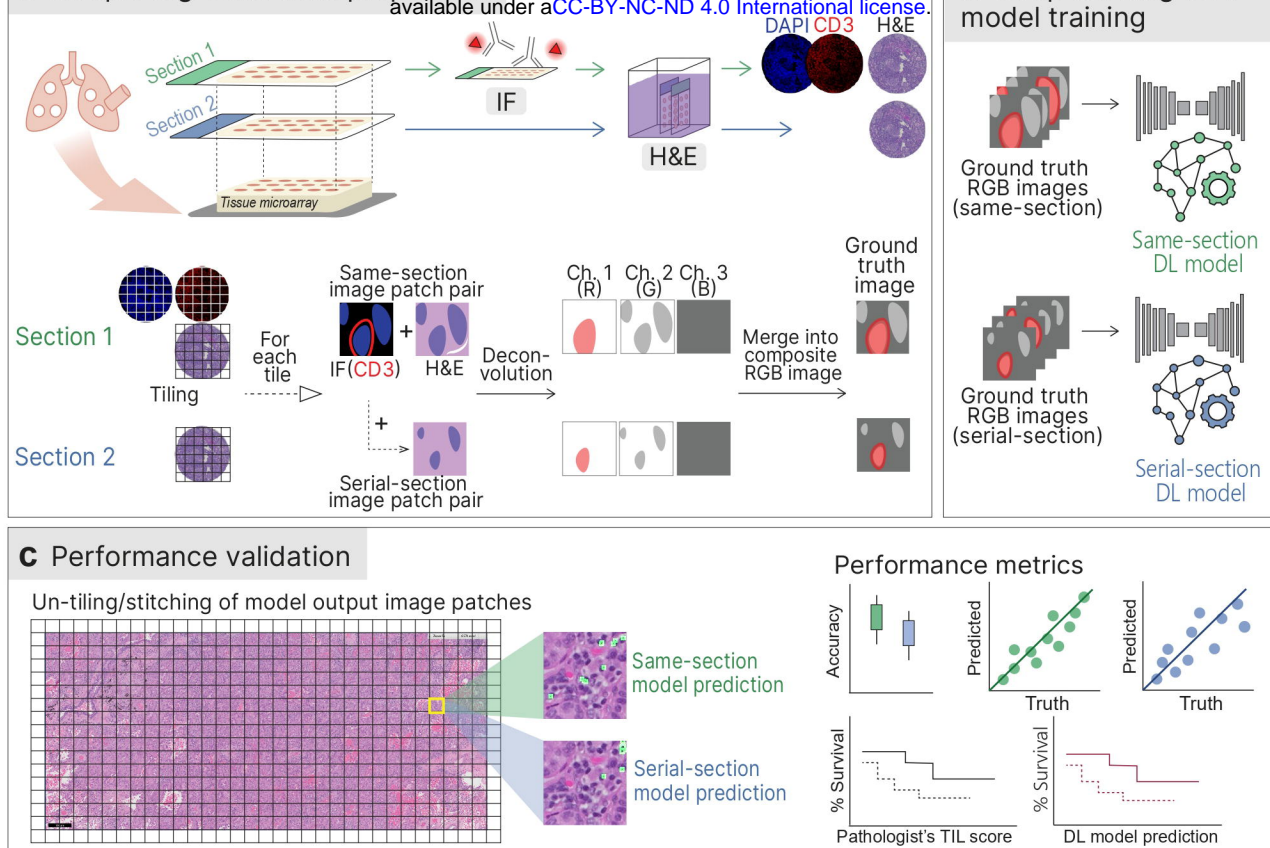410    Available from https://doi.org/10.48550/arXiv.1611.07004.

19

Figure 1: Schematic diagram of the study protocol. a) Preparation of samples and ground truth for both serial-section and same-section datasets; b) Construction and training of two P2P-GAN DL models utilizing the serial-section and same-section datasets; c)  Validation of DL model performance using an independent in-house IHC cohort and an external lung cancer cohort.
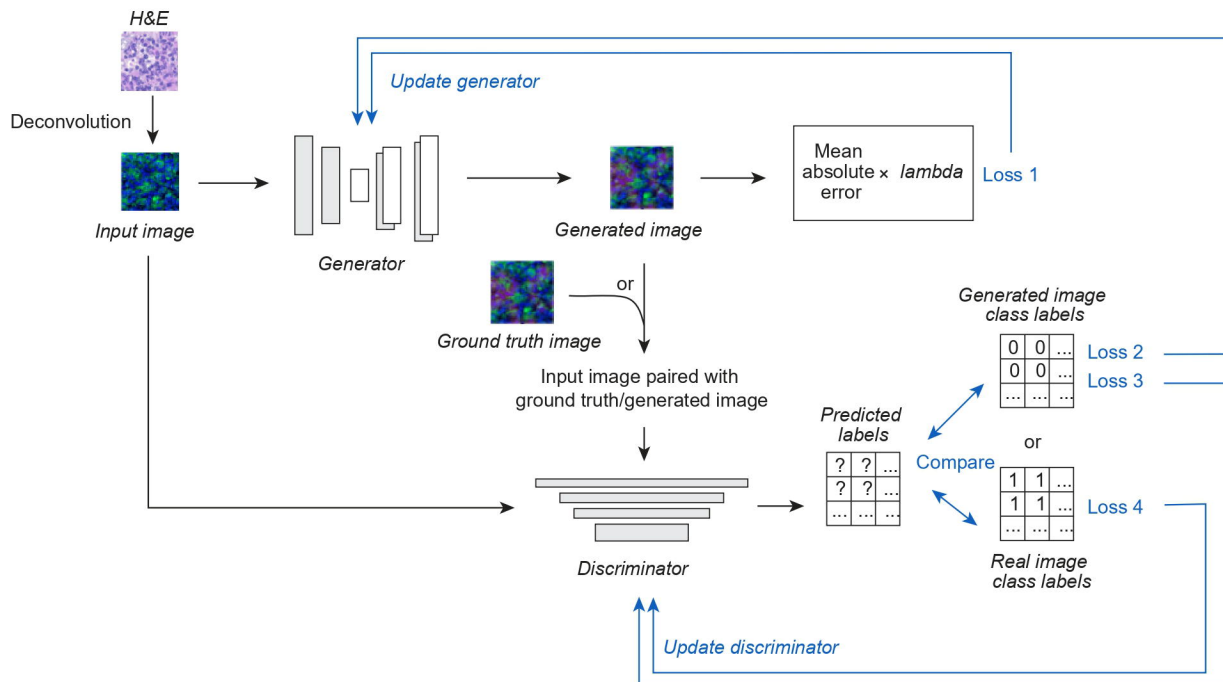
Figure 2: P2P-GAN model architecture and parameter updating process during model training. Two key components, namely the generator, which inputs the H&E image patch and generates (predicts) CD3+ signals on the input image, and its adversary, the discriminator, which distinguishes the generator output from the image with true CD3+ signals (ground truth). The adversarial nature of the network enables the generator to produce good predictions.
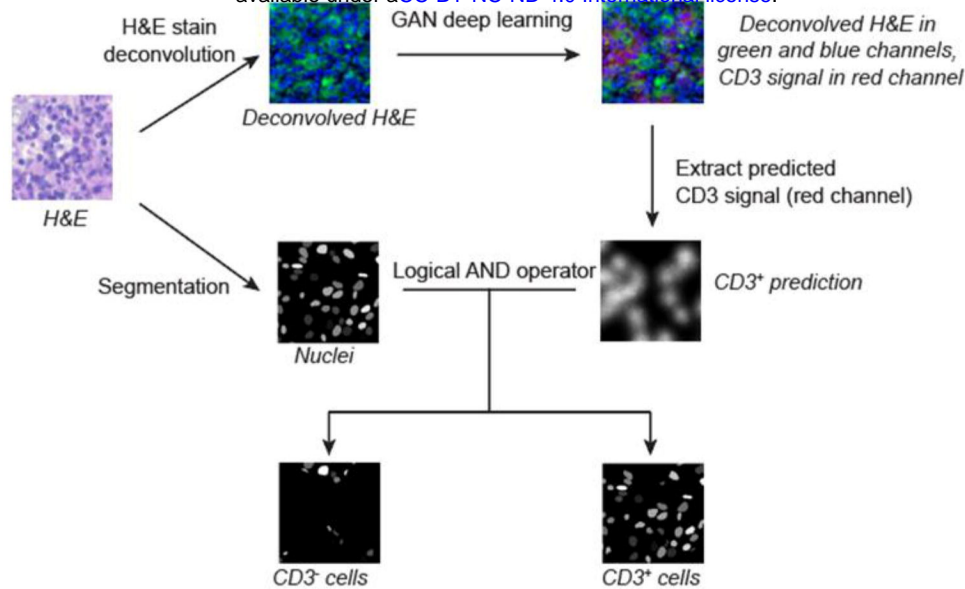
Figure 3: Identification of CD3$^+$ and CD3$^-$ T-cells predicted by our proposed P2P-GAN models. The process involves extracting the model-predicted CD3 signals (in the red channel) and overlaying the detected signal onto nuclei detected in the H&E image. Nuclei with matching CD3 signals are regarded as CD3$^+$ T-cells, otherwise the nuclei are regarded as CD3$^-$ T-cells.
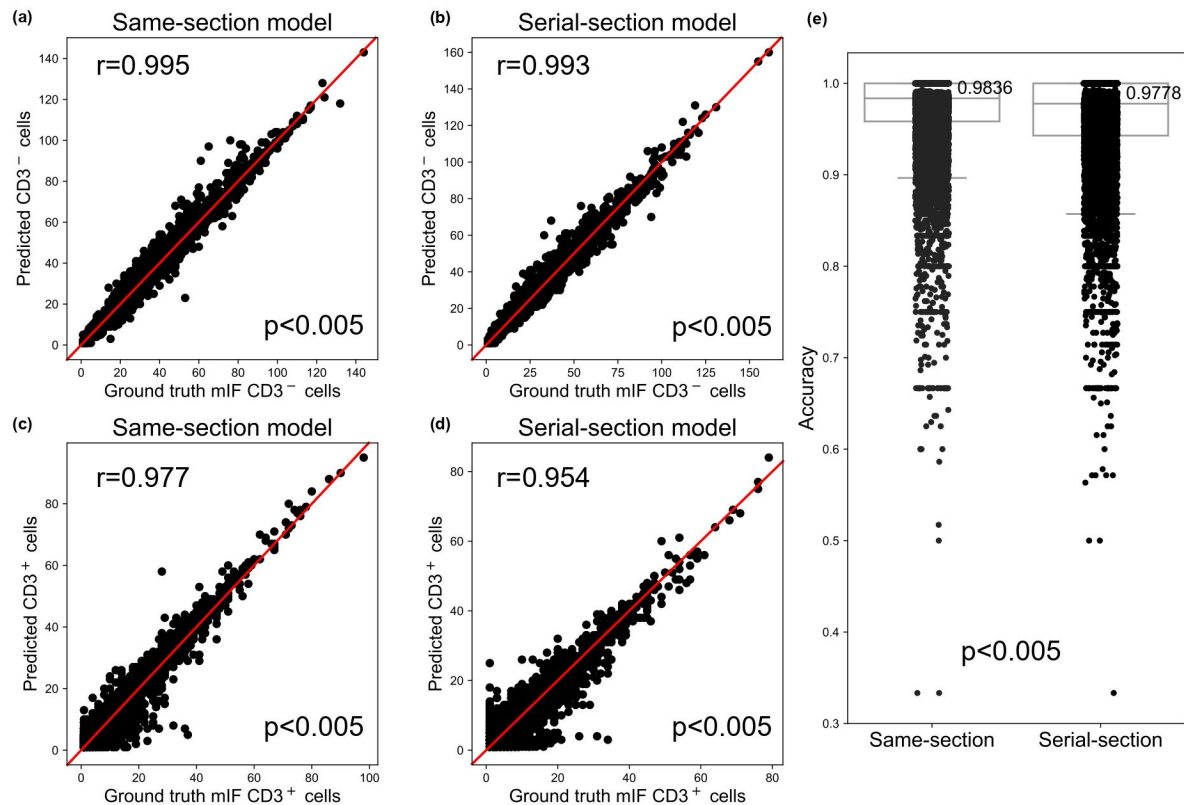
Figure 4: Model performance evaluation using the corresponding training cohorts (i.e., same-section and serial-section datasets, respectively). Comparison of model-predicted (a-b) CD3$^-$ and (c-d) CD3$^+$ cell (y-axis) counts with mIF-quantified CD3$^+$ cell counts (x-axis) using Pearson's correlation analysis. (e) Overall accuracy comparison between the model prediction accuracy (y-axis) of the same-section (left) and serial-section (right) models, using the randomly selected held-out samples from the same-section training cohort based on Mann-Whitney U-tests; each dot represents an image patch.
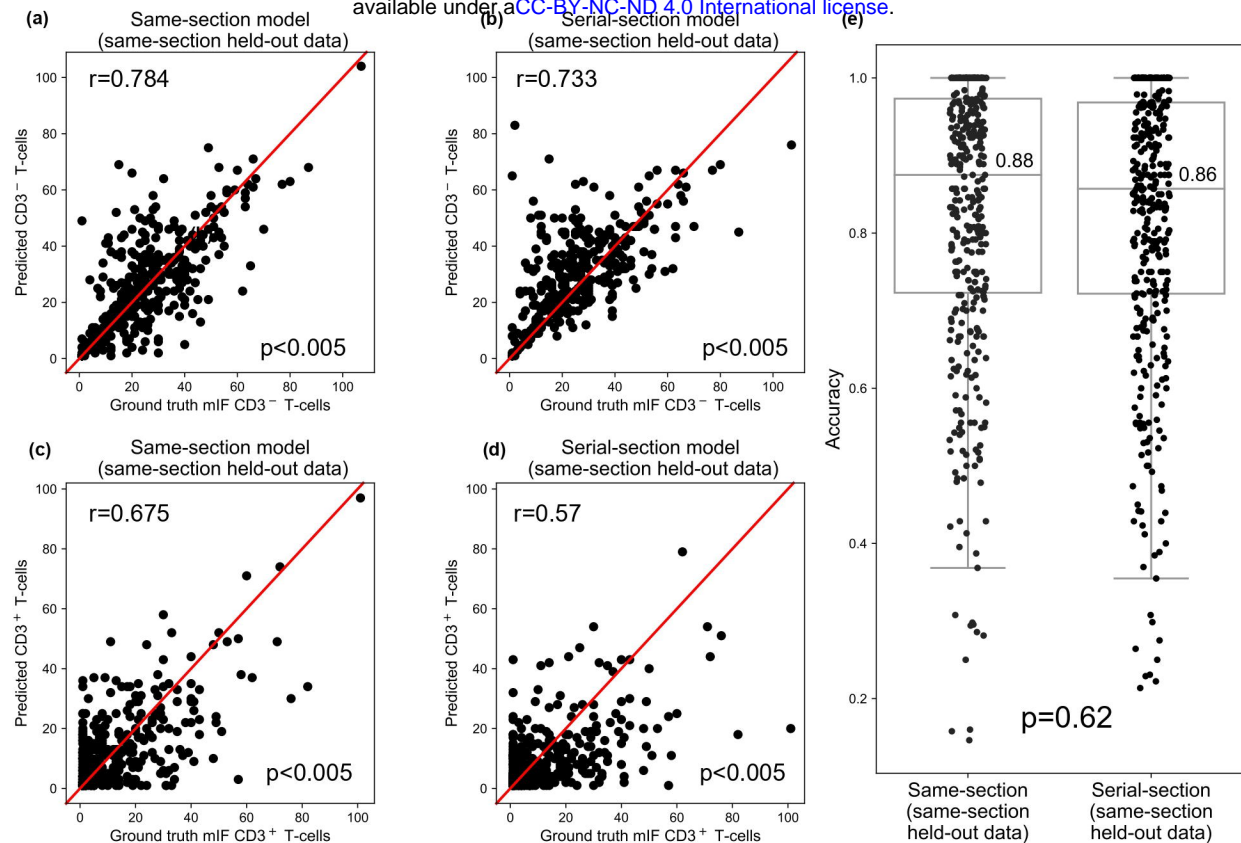
Figure 5: Model performance evaluation using the randomly selected held-out samples from the same-section training cohort. Comparison of model-predicted (a-b) CD3⁻ and (c-d) CD3⁺ cell (y-axis) counts with mIF-quantified CD3⁺ cell counts (x-axis) using Pearson's correlation analysis. (e) Overall accuracy comparison between the model prediction accuracy (y-axis) of the same-section (left) and serial-section (right) models, using the randomly selected held-out samples from the same-section training cohort based on Mann-Whitney U-tests; each dot represents an image patch.
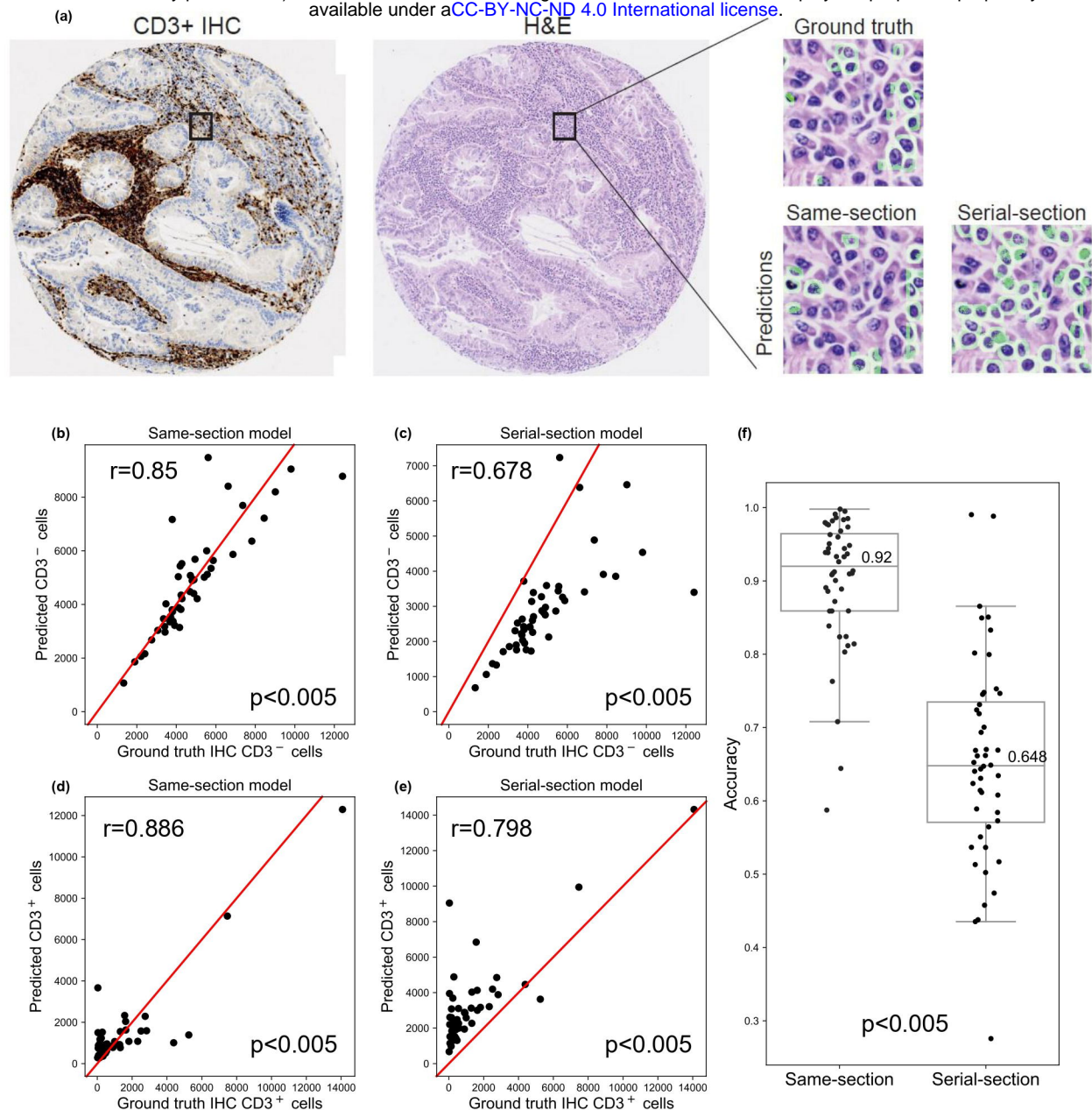
Figure 6: Model performance evaluation using the randomly selected held-out samples from the same-section training cohort with representative images along with model predicted CD3$^+$ T-cell visualization in (a). Comparison of model-predicted (b-c) CD3$^-$ and (d-e) CD3$^+$ (y-axis) cell counts with mIF-quantified CD3$^+$ cell counts (x-axis) using Pearson's correlation. (f) Overall accuracy comparison between the model prediction accuracy (y-axis) of the same-section (left) and serial-section (right) models, using the randomly selected held-out samples

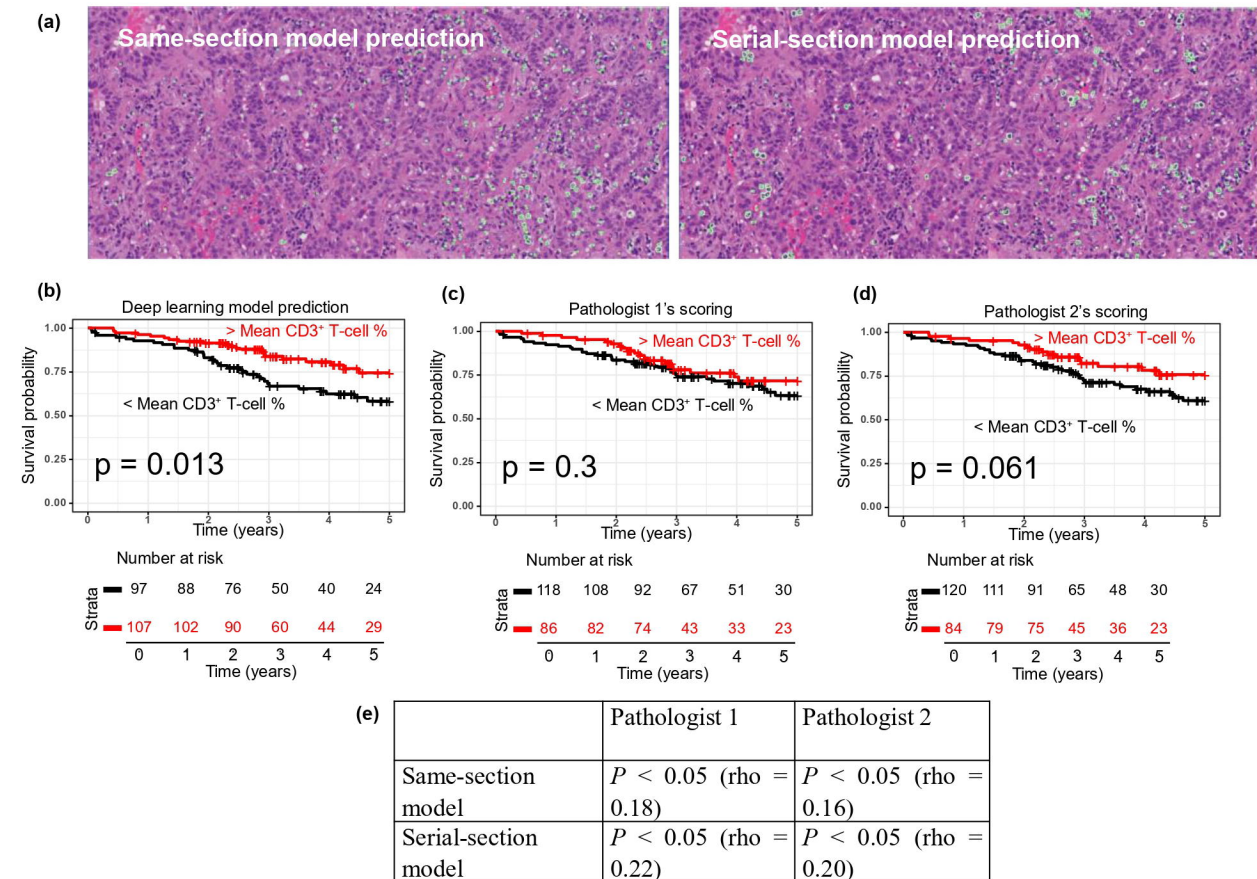from the same-section training cohort based on Mann-Whitney U-tests; each dot represents an image patch.



(a)

(b) Deep learning model prediction
Survival probability
> Mean CD3+ T-cell %
< Mean CD3+ T-cell %
p = 0.013
Time (years)

Number at risk
Strata
97 88 76 50 40 24
107 102 90 60 44 29
0 1 2 3 4 5
Time (years)

(c) Pathologist 1's scoring
Survival probability
> Mean CD3+ T-cell %
< Mean CD3+ T-cell %
p = 0.3
Time (years)

Number at risk
Strata
118 108 92 67 51 30
86 82 74 43 33 23
0 1 2 3 4 5
Time (years)

(d) Pathologist 2's scoring
Survival probability
> Mean CD3+ T-cell %
< Mean CD3+ T-cell %
p = 0.061
Time (years)

Number at risk
Strata
120 111 91 65 48 30
84 79 75 45 36 23
0 1 2 3 4 5
Time (years)

(e)

|  | Pathologist 1 | Pathologist 2 |
|---|---|---|
| Same-section model | $P < 0.05$ (rho = 0.18) | $P < 0.05$ (rho = 0.16) |
| Serial-section model | $P < 0.05$ (rho = 0.22) | $P < 0.05$ (rho = 0.20) |

Figure 7: Model performance evaluation using an external lung cohort, with representative images along with model predicted CD3$^+$ T-cell visualization in (a). Survival analyses using the external lung cohort show (b) significant association between (same-section) model-predicted CD3 patient groups (low versus high %CD3$^+$ abundance groups based on the average %CD3$^+$ T-cell counts), while no significant association was observed with the use of manual TIL scoring by (c) pathologist 1, and (d) pathologist 2. (e) Spearman correlation of the prediction of CD3$^+$ densities using the same- and serial-section models with the manual TIL density scoring by two independent pathologists.

Table 1: Cohort characteristics

| Dataset type | Number of patients | Images per patient | Image modalities | Tissue format | Image size (pixels) |
|---|---|---|---|---|---|
| mIF training dataset (same- and serial-section) | 57 for both the same-section and serial-section datasets | 1 | H&E and mIF | TMA cores | 3228´3228 |
| IHC testing dataset (serial-section) | 48 | 1 | H&E and IHC | TMA cores | Approximately 4000´4000 |
| Onco-SG testing dataset[18] | 204 | 1–3 | H&E | Region of interest in resected tissues | 1792´768 |