

Deep learning detects virus presence in cancer histology

Jakob Nikolas Kather^{1,2,3,4}, Jefree Schulte⁵, Heike I. Grabsch^{6,7}, Chiara Loeffler¹, Hannah Muti¹,
James Dolezal⁴, Andrew Srisuwananukorn⁸, Nishant Agrawal⁹, Sara Kochanny⁴, Saskia von Stillfried¹⁰,
Peter Boor¹⁰, Takaki Yoshikawa^{11,12}, Dirk Jaeger^{2,3}, Christian Trautwein¹, Peter Bankhead¹³,
Nicole A. Cipriani⁵, Tom Luedde¹, Alexander T. Pearson⁴

¹ Department of Medicine III, University Hospital RWTH Aachen, Aachen, Germany

² German Cancer Consortium (DKTK), German Cancer Research Center (DKFZ), Heidelberg, Germany

³ Applied Tumor Immunity, German Cancer Research Center (DKFZ), Heidelberg, Germany

⁴ Department of Medicine, University of Chicago, Chicago, IL, USA

⁵ Department of Pathology, University of Chicago Medicine, Chicago, IL, USA

⁶ Pathology & Data Analytics, Leeds Institute of Medical Research at St James's, University of Leeds, Leeds,
UK

⁷ Department of Pathology, GROW School for Oncology and Developmental Biology, Maastricht University
Medical Center+, Maastricht, The Netherlands

⁸ Department of Medicine, University of Illinois – Chicago, Chicago, IL, USA

⁹ Department of Surgery, University of Chicago, Chicago, IL, USA

¹⁰ Department of Pathology, University Hospital RWTH Aachen, Aachen, Germany

¹¹ Department of Gastrointestinal Surgery, Kanagawa Cancer Center, Yokohama, Japan

¹² Department of Gastric Surgery, National Cancer Center Hospital, Tokyo, Japan

¹³ MRC Institute of Genetics and Molecular Medicine, University of Edinburgh, Edinburgh, UK

1 **Abstract**

2 Oncogenic viruses like human papilloma virus (HPV) or Epstein Barr virus (EBV) are a major cause of human
3 cancer. Viral oncogenesis has a direct impact on treatment decisions because virus-associated tumors can
4 demand a lower intensity of chemotherapy and radiation or can be more susceptible to immune check-
5 point inhibition. However, molecular tests for HPV and EBV are not ubiquitously available.

6 We hypothesized that the histopathological features of virus-driven and non-virus driven cancers are suf-
7 ficiently different to be detectable by artificial intelligence (AI) through deep learning-based analysis of
8 images from routine hematoxylin and eosin (HE) stained slides. We show that deep transfer learning can
9 predict presence of HPV in head and neck cancer with a patient-level 3-fold cross validated area-under-
10 the-curve (AUC) of 0.89 [0.82; 0.94]. The same workflow was used for Epstein-Barr virus (EBV) driven
11 gastric cancer achieving a cross-validated AUC of 0.80 [0.70; 0.92] and a similar performance in external
12 validation sets. Reverse-engineering our deep neural networks, we show that the key morphological fea-
13 tures can be made understandable to humans.

14 This workflow could enable a fast and low-cost method to identify virus-induced cancer in clinical trials or
15 clinical routine. At the same time, our approach for feature visualization allows pathologists to look into
16 the black box of deep learning, enabling them to check the plausibility of computer-based image classifi-
17 cation.

18

19 **Introduction**

20 Oncogenic viruses cause approximately 15% of malignant tumors in humans.¹ Viruses can induce cancers
21 with different histology and across different anatomic sites including squamous cell carcinomas (e.g. head
22 and neck, cervix), adenocarcinomas (e.g. gastric), sarcomas (e.g. Kaposi), lymphomas (e.g. Burkitt) and
23 hepatocellular carcinoma. Virus-driven tumors are an important health issue in western countries, but
24 their global health impact is even higher as 80% of all virus-driven cancers occur in developing nations.²
25 Their incidence is expected to increase drastically in the next decade in developing and economically de-
26 veloped countries.^{3,4} Some types of cancer are almost always virally driven (e.g. cervical cancer) while
27 others can have viral or non-viral driver mechanisms (e.g. head and neck cancer or gastric cancer). In these
28 cases, it is important to determine if a patient's tumor has a viral origin because if this is the case, a dif-
29 ferent clinical management may be warranted and virus status might influence the choice of a clinical trial
30 for that particular patient. For example, in the case of head and neck squamous cell carcinoma (HNSC),
31 patients with human papilloma virus (HPV)-positive tumors have superior overall survival compared to
32 patients with HPV-negative tumors of the same stage and can benefit from treatment de-escalation.⁵
33 Likewise, patients with Epstein-Barr-Virus (EBV) related gastric adenocarcinoma tend to have a better
34 prognosis and EBV positivity has been suggested as a biomarker for immunotherapy response.⁶

35 The gold standard method for detection of viruses in human cancer is dependent on the tumor type. In
36 head and neck cancer, overexpression of p16 as assessed by immunohistochemistry is the most commonly
37 used surrogate marker for virus presence. However, p16 is neither perfectly sensitive nor specific⁷, and
38 some centers also use HPV polymerase-chain reaction, in-situ hybridization, or targeted DNA sequencing
39 for HPV detection in tumor tissue. While these tests are more specific, they are also more expensive and
40 time consuming. Presence of latent EBV infection in gastric cancer is usually measured using EBV-encoded
41 RNA in-situ hybridization in pathology samples, which has a relatively high sensitivity and specificity but
42 requires dedicated testing equipment and expertise for accurate interpretation.

43 In the present study, we hypothesized that morphological features correlating with the presence of vi-
44 ruses in solid tumors can be deduced from hematoxylin and eosin (H&E) histology, which is routinely
45 available for almost any patient with a solid tumor. As a tool for feature extraction from images, we used
46 deep learning, a form of artificial intelligence (AI), which has previously been used to detect high-level
47 morphological features directly from histological images.⁸⁻¹⁰

48

49 **Methods**

50 **Ethics and data sources**

51 All experiments were conducted in accordance with the Declaration of Helsinki and the International Eth-
52 ical Guidelines for Biomedical Research Involving Human Subjects. Anonymized scanned whole slide im-
53 ages were retrieved from The Cancer Genome Atlas (TCGA) project through the Genomics Data Commons
54 Portal (<https://portal.gdc.cancer.gov/>). From this source, we retrieved images of head and neck squamous
55 cell carcinoma (HNSC)¹¹ and gastric adenocarcinoma (stomach adenocarcinoma, STAD)¹². Exclusion crite-
56 ria for patients in these cohorts were missing values in virus status, corrupt image files or lack of tumor
57 tissue on the whole slide image. For TCGA-HNSC, images from N=450 patients were downloaded of which
58 N=38 met exclusion criteria, leaving images from N=412 patients for further processing. For TCGA-STAD,
59 images from N=416 patients were downloaded of which N=99 met exclusion criteria, leaving N=317 pa-
60 tients for further processing. Furthermore, we retrieved anonymized archival tissue samples of N=105
61 patients with HNSC from the University of Chicago Medicine Pathology archive (Chicago, Illinois, USA;
62 “UCH-HNSC”) and anonymized tissue samples of N=197 patients with gastric cancer from the Kanagawa
63 Cancer Center Hospital (Yokohama, Japan; “KCCH-STAD”) as described before¹³. For HNSC, HPV status was
64 determined as described by Campbell et al.¹⁴ (by consensus of DNA sequencing¹⁵ and RNA sequencing¹⁶).
65 For TCGA-STAD, EBV status was retrieved from genomic subtypes as described by Liu et al.¹⁷. For samples
66 in UCH-HNSC, HPV status was defined by polymerase-chain reaction for the viral genes E6 and E7. For
67 tumor samples in KCCH-STAD, EBV status was defined by EBV-encoded RNA in-situ hybridization.¹⁸

68 **Deep transfer learning workflow**

69 All histological slides were reviewed and tumor regions were manually delineated in QuPath¹⁹, tessellated
70 into tiles of 256 x 256 μm^2 which were subsequently downsampled to 224 x 224 px, yielding an effective
71 magnification of 1.14 $\mu\text{m}/\text{px}$. These tumor tiles were used for deep transfer learning in MATLAB R2019a

72 as described before^{9,10}. We used a modified VGG19 deep convolutional neural network²⁰ which was pre-
73 trained on ImageNet (<http://www.image-net.org>, architecture shown in Suppl. Table 1). VGG19 was cho-
74 sen because of its previously proven performance in detecting multiple tissue components in human can-
75 cer histology⁹ and because of its compatibility with the Deep Dream method (see below). All TCGA cohorts
76 were randomly split into three equal subsets at patient level. A VGG19 classifier was trained on these data
77 in a three-fold cross-validated way. This procedure yielded three independent classifiers which were eval-
78 uated on their respective test set of held-out patients. For each tumor type, the classifier was subse-
79 quently re-trained on the whole TCGA set and evaluated on an external test set.

80 **Feature visualization**

81 To trace back deep-learning based predictions to human-understandable morphological patterns in his-
82 tology, we used deep-dream-based visualization of output layer neurons for each class. We used a
83 MATLAB implementation (<https://de.mathworks.com/help/deeplearning/ref/deepdreamimage.html>) of
84 the original DeepDream algorithm ([https://github.com/tensorflow/tensorflow/blob/master/tensor-](https://github.com/tensorflow/tensorflow/blob/master/tensorflow/examples/tutorials/deepdream/deepdream.ipynb)
85 [flow/examples/tutorials/deepdream/deepdream.ipynb](https://github.com/tensorflow/tensorflow/blob/master/tensorflow/examples/tutorials/deepdream/deepdream.ipynb)) with pyramid level 6 and 500 iterations and sub-
86 sequently auto-optimized colors by histogram stretching in IrfanView 4.52 (<https://www.irfanview.com/>).
87 Color optimization was done with identical parameters for all deep-dream-images generated by a given
88 network.

89 **Statistics and data presentation**

90 Classifier performance was assessed by the Area under the Receiver Operating Curve (AUC under the ROC)
91 with sensitivity (true positive rate, TPR) plotted on the vertical axis and 1 – specificity (false positive rate,
92 FPR) plotted on the horizontal axis. 95% confidence intervals for the AUC were calculated with 500-fold
93 bootstrapping with the “bias corrected and accelerated percentile method”²⁸ unless otherwise stated.
94 For three-fold cross validated experiments, the mean of AUCs and the mean of confidence interval

95 boundaries from all three classifiers is given if not otherwise noted. The ROC procedure is a widely used
96 technique to assess the power of a classifier for any possible cutoff value of a numerical test. In this study,
97 the cutoff for “percentage of virus-positive image tiles” was varied, yielding different sensitivity/specificity
98 pairs which are plotted as ROC curves.

99 **Data availability**

100 Images from the TCGA cohorts are available at <https://portal.gdc.cancer.gov/>. Our source codes are avail-
101 able at <https://github.com/jnkather/VirusFromHE>.

102

103 **Results**

104 **Deep learning detects virus presence in squamous cell carcinomas and adenocarcinomas**

105 We hypothesized that the presence of human papillomavirus (HPV) can be detected in head and neck
106 squamous cell carcinoma (HNSC, Figure 1a) and that the presence Epstein-Barr-Virus (EBV) can be de-
107 tected in gastric adenocarcinoma (STAD, Figure 1b) directly from histology by deep learning with a convo-
108 lutional neural network (CNN). We used hematoxylin and eosin (H&E) stained tissue slides of patients
109 included in the multicenter TCGA cohort (Suppl. Table 2) and trained a deep learning classifier in a patient-
110 level three-fold cross-validated way (Figure 1c), followed by re-training on the whole cohort (Figure 1d).
111 In head and neck cancer (N=412 patients, 12% HPV positive), this yielded an average patient-level AUC of
112 0.89 [0.82; 0.94] (Figure 1e) and applying the same workflow to detect EBV in gastric cancer (N=317 pa-
113 tients from TCGA, 8% EBV positive, Suppl. Table 3), a patient-level three-fold cross-validated neural net-
114 work achieved an AUC of 0.80 [0.70; 0.92] (Figure 1f). Together, these results show that deep learning can
115 robustly distinguish virus-induced (“virus present”) from non-virus-induced tumors (“virus not present”)
116 across different histologies and anatomic sites.

117 **Noisy tile level data yields high patient-level accuracy in external validation cohorts**

118 To assess the robustness of the classifiers, we used the neural network that was trained on the entire
119 TCGA patient cohorts for head and neck and gastric cancer, respectively, and evaluated the classifiers on
120 external validation cohorts. Non-overlapping tissue tiles of 256 μm edge length were used to predict a
121 “virus probability score” which classified each tile as either virus positive or negative (derived from a tu-
122 mor that was virally induced or derived from a tumor that was non-virally induced). These predictions
123 were subsequently pooled on a patient level as “fraction of positive tiles” with varying thresholds accord-
124 ing to the Receiver Operating Characteristic procedure (Figure 2a). Because each tile in the training set
125 was assigned the label of the corresponding patient (obtained via bulk testing of tissue) and the tiles

126 contained a multitude of different tissue types (tumor epithelium, stroma, necrosis, mucus, and others),
127 this training set was inherently noisy. Correspondingly, predictions for virus-negative tiles in the EBV test-
128 ing set were noisy with many false positive tile-level predictions (Figure 2b, right-hand side). However,
129 tiles from virus-positive patients were mostly classified correctly (Figure 2a, left panel), enabling robust
130 classification after pooling tile-level predictions on a patient level. AUC for EBV detection in the KCCH-
131 STAD cohort (N=197 patients, 5% EBV positive) was 0.81 [0.69; 0.89] (Figure 2c; trained on TCGA-STAD,
132 tested on KCCH-STAD). We manually assessed the histomorphology of tissue tiles in the KCCH-STAD co-
133 hort (Figure 2d, more examples are available at <http://doi.org/10.5281/zenodo.3247009>) and found that
134 false positive tiles often presented with lymphocyte-rich stroma, a known morphological hallmark of EBV-
135 positive gastric cancer.²¹ Thus, we conclude that misclassifications of individual image tiles were due to
136 plausible human-understandable morphological features.

137 Similarly, we validated the virus detector for HPV trained on TCGA-HNSC (N=412 patients, 12% HPV posi-
138 tive) in our in-house cohort UCH-HNSC from University of Chicago (N=105 patients, 49% HPV positive).
139 This cohort had two main differences compared to the TCGA cohort which might negatively affect classi-
140 fier performance: first, a polymerase-chain reaction for high-risk HPV viral genes was used to determine
141 virus status. Second, this cohort was artificially balanced for HPV status and thus had a much higher prev-
142 alence of HPV-induced cancer than TCGA. In spite of these stark differences, our classifier achieved an
143 AUC of 0.70 [0.66; 0.74] for HPV prediction in UCH-HNSC. Manual review of representative tissue tiles by
144 an expert pathologist showed that tiles with a high HPV prediction score were “carcinomas with large
145 nested, broad-based invasion and relative decrease in cytoplasmic keratinization, resulting in a blue (cool-
146 toned) appearance”. This is compatible with previously known morphological features of HPV-positive
147 HNSC.²² Tiles with a low HPV prediction score were “carcinomas with small nested invasion and eosino-
148 philic (pink or warm-toned) cytoplasmic keratinization”. Thus, we conclude that in HNSC as well as in

149 gastric cancer, predictions of viral status in individual tissue tiles by a deep neural network were plausible
150 to expert pathologists.

151 Together, these data show that despite noisy training data and tile-level misclassifications, patient-level
152 prediction of virus status in HNSC and gastric cancer can reach a high accuracy.

153 **Reverse-engineering trained neural networks**

154 Attempting to characterize more precisely which morphological features may have been used by the neu-
155 ronal network to detect virus-induced cancer, we used a feature-visualization method and discussed the
156 results with a panel of expert pathologists. We hypothesized that reverse-engineering features from neu-
157 ral networks could be used as a plausibility check for deep learning, completing the cycle “human to AI
158 and back”. In an exploratory study, we employed the Deep Dream algorithm which uses a trained neural
159 network (Figure 3a) to create pseudo-images for each output class in the classification layer (Figure 3b).
160 This approach yielded “pseudo-histology” images for HPV positive and negative HNSC (Figure 3b) and EBV
161 positive and negative gastric cancer (Figure 3c), discussing the resulting images with five pathologists. In
162 general, pathologists described the images as “beautiful” and “psychedelic”. Relating the aspect of
163 pseudo-histology in histological terms, they described the features as “a sheet of small nodules composed
164 of bright, predominantly warm colors” (HPV negative HNSC, Figure 3c, left panel), “large nests with
165 rounded borders composed of dark, predominantly cool colors punctuated by red dots” (HPV positive
166 HNSC, Figure 3c, right panel) and “ill-defined dark green whorls punctuated by blue dots and wisps of
167 yellow” (EBV negative gastric cancer [STAD], Figure 3d, left panel), “overlapping sheets with reticulated
168 patterns in pastel colors”, potentially resembling “lymphoid stroma” (EBV positive gastric cancer, Figure
169 3d right panel).

170 Together, these data show that deep learning can plausibly sort tissue tiles (Figure 2d) and yields a high
171 classification performance for virus presence (Figure 2c). The actual morphological patterns used for this

172 classification may be different from the ones that humans typically use but can be visualized in a way that
173 might be understandable for humans through the Deep Dream algorithm (as has been shown in non-
174 medical applications²³). Based on this, we conclude that, by analyzing tile-level classification and possibly
175 by analyzing Deep Dream images, human observers can get an insight about morphological patterns used
176 by deep neural networks, allowing for quality control and possibly performing human classification per-
177 formance through machine-identified morphological features.

178

179 **Discussion**

180 Virus-induced cancers mostly occur in developing countries, making them neglected diseases on a global
181 scale. Some of these virus-driven cancers are under-tested for in clinical routine. Existing wet lab assays
182 to test for virus presence (such as sequencing) are costly, require a high level of expertise (such as in-situ
183 hybridization) and not all assays achieve a perfect classification accuracy (such as p16 immunohistochem-
184 istry²⁴). Here, we present a deep-learning-based low-to-no-cost assay for routine detection of virus pres-
185 ence from ubiquitously available histology in two major tumor types of very different histology. We
186 demonstrate that classification accuracy is as high as AUC 0.81 when trained with a few hundred patients.
187 Our approach relies on digitally scanned images of hematoxylin & eosin stained tissue slides. The cost to
188 scan such a histology slide is well below \$10 at low throughput and considerably lower at high through-
189 put¹⁰, potentially enabling noticeable cost savings for virus testing of tumor tissue in the future.

190 At the moment, sensitivity and specificity of our classifier is lower than in routine diagnostic tests: for EBV
191 detection in gastric cancer by EBV-encoded RNA in-situ hybridization, one study reported a sensitivity of
192 100% at a specificity of 90%.²⁵ For HPV detection in HNSC by p16 immunohistochemistry, another study
193 reported a sensitivity of 97.4% and a specificity 93.75%.²⁴ As shown in Figure 1e-f and Figure 2c, the deep
194 learning classifier approaches these gold standard methods but is still less sensitive and specific. However,
195 sensitivity and specificity of our method are higher than those in previous studies of deep-learning based
196 prediction of molecular features from histology.^{8,26}

197 In our experiments, the deep learning classifiers reached a high cross-validated performance which we
198 could replicate in an external validation set for gastric cancer (Figure 2c). In the multicenter TCGA-HNSC
199 cohort, cross-validated HPV detection performance was high, but dropped in the external validation co-
200 hort UCH-HNSC. This may be related to the relatively small patient size of this cohort or due to different
201 gold standards for HPV detection (consensus of DNA and RNA sequencing in TCGA-HNSC and polymerase-

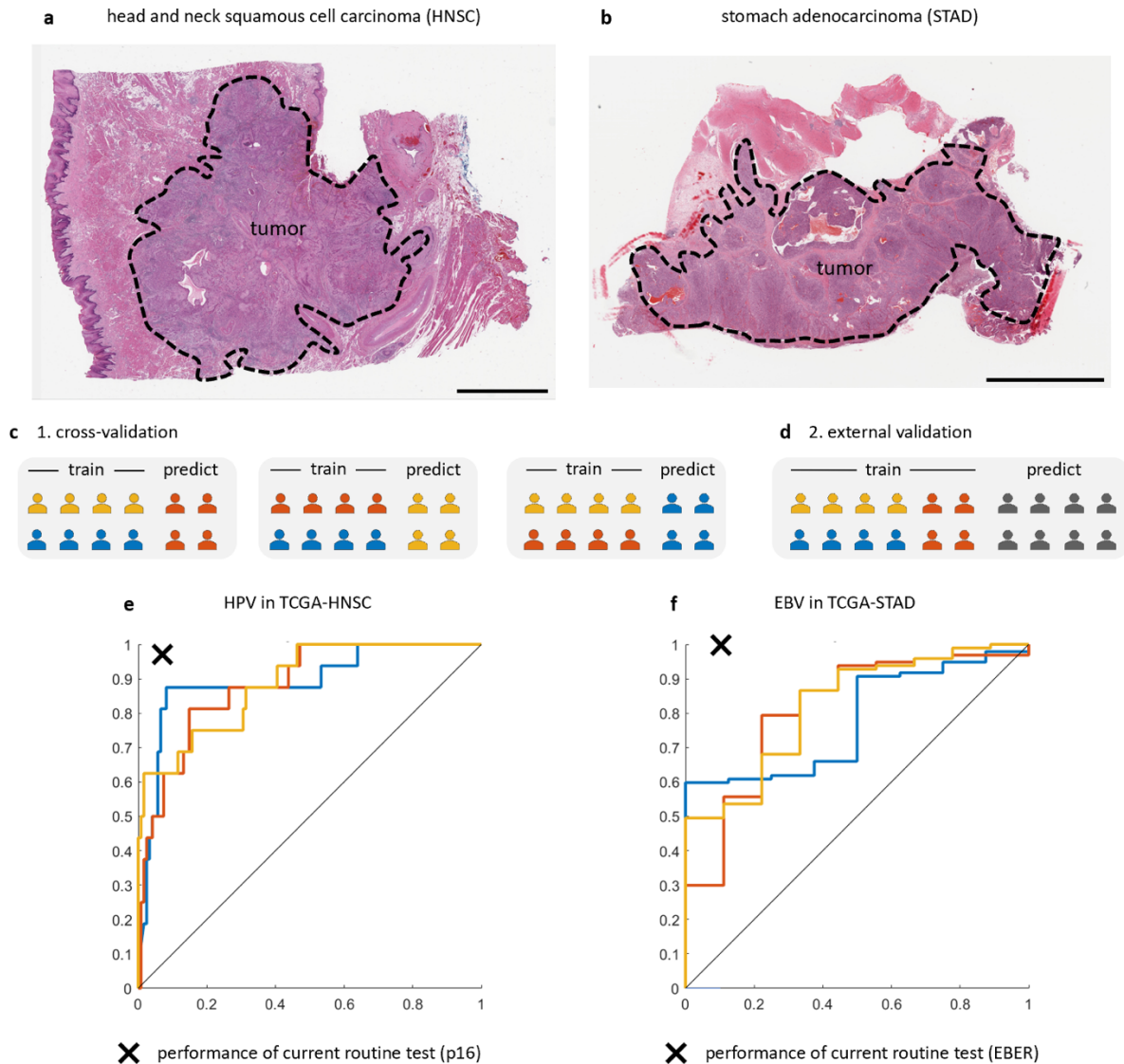
202 chain reaction in UCH-HNSC). Most probably, however, this is due to the very different prevalence of virus-
203 induced cancers in the training set and in the test set. Whereas the training set (TCGA-HNSC) reflected
204 the natural prevalence of HPV-positive cancers, the test set (UCH-HNSC) was artificially balanced to a
205 prevalence of 50%. This may have negatively affected classifier performance as has been described for
206 mutation prediction in lung cancer⁸.

207 According to our experience from similar tasks, it can be expected that training on larger clinical cohorts
208 will likely improve performance of our method. Similarly, further optimizing hyperparameters and neural
209 network architectures will likely yield a performance boost. Also, further dividing deep learning classifiers
210 by anatomical sub-sites of tumors (e.g. oropharyngeal or hypopharyngeal) will likely increase perfor-
211 mance. In the end, this image-based biomarker, like all biomarkers, needs to be tested in prospective
212 clinical trials before widespread clinical use.

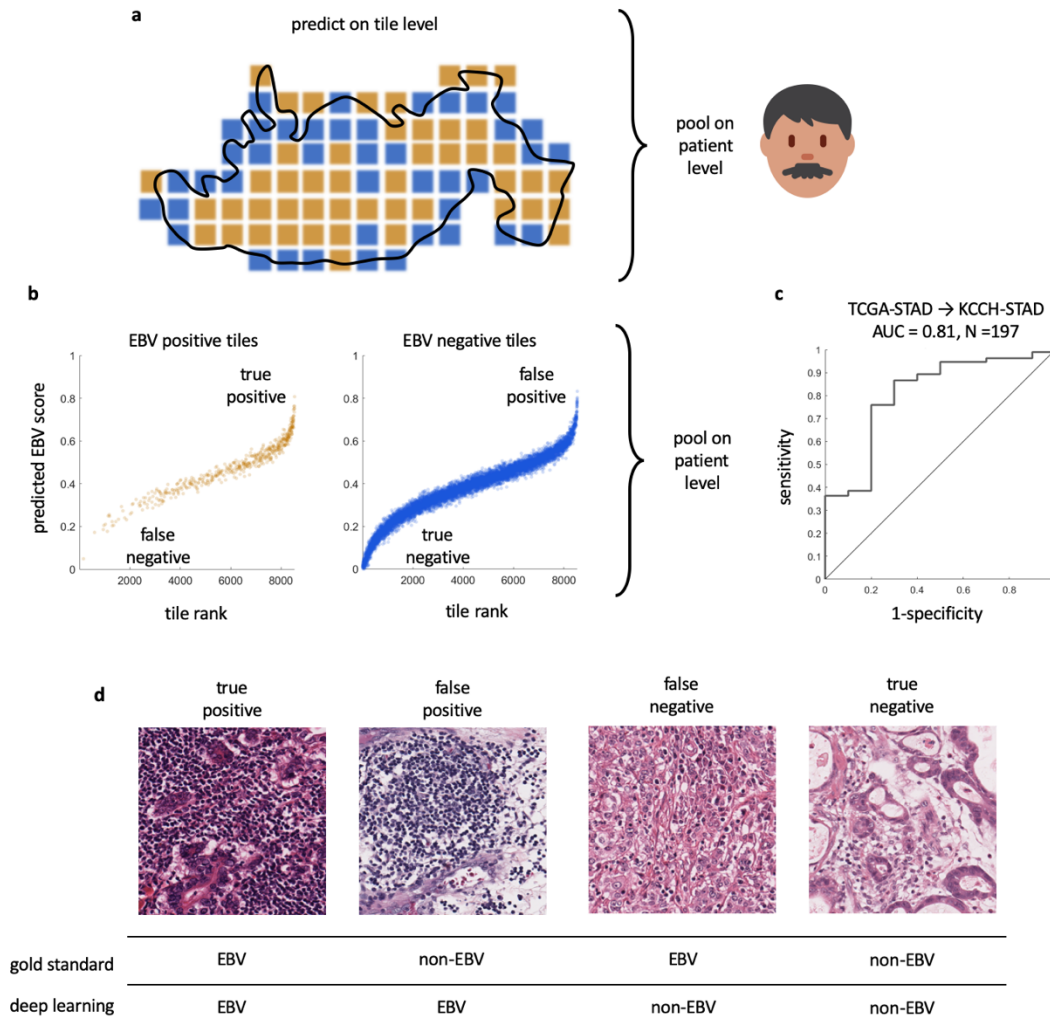
213 A new aspect of our study is the approach “human to AI and back”: humans (expert pathologists) deline-
214 ated tumor tissue in whole slide sections and thus enabled the AI to detect virus presence in histological
215 images. In turn, using deep-dream-based feature visualization, we show that the AI can in principle inform
216 a human observer about morphological features of interest. Feature visualization by Deep Dream and
217 similar methods²⁷ is well-established to understand the inner workings of deep neural networks. Yet, to
218 our knowledge, this has never been systematically used for pathologist-AI-crosstalk. Thus, our study
219 shows for the first time that deep learning algorithms can not only be used as tools to facilitate diagnostic
220 routine but could also enable human observers to get a different viewpoint on histomorphology.

221

222 Figures

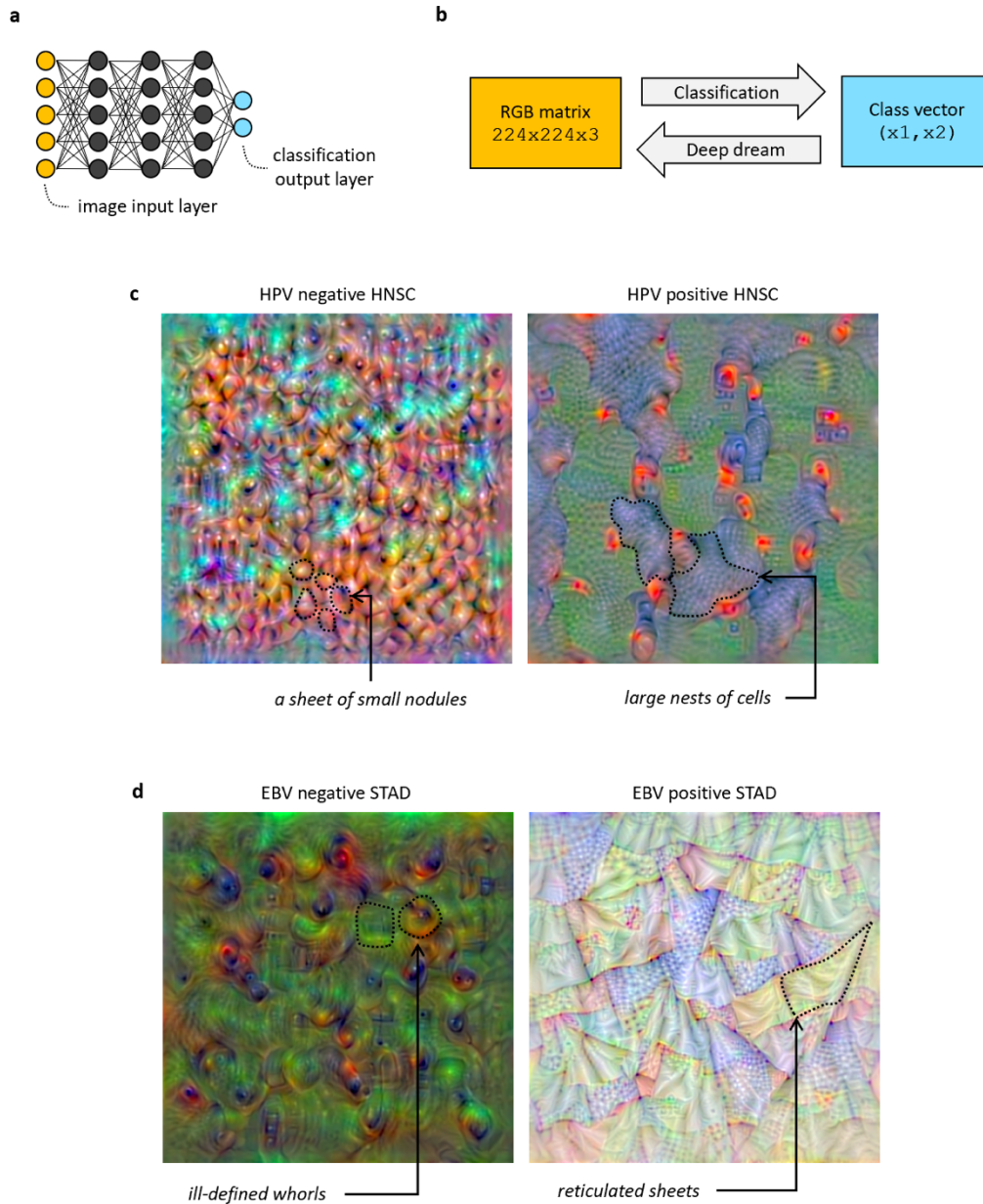


223 **Figure 1: Detection of virus-induced cancer from histological routine images.** (a) A representa-
224 tive sample (one of N=412 patients) of the TCGA-HNSC training cohort with a manual outline
225 around the tumor tissue, (b) a representative sample of the TCGA-STAD cohort (1 of N=317 pa-
226 tients), scale bars in a and b are 5000 μ m. (c) We trained and tested classifiers with patient-level
227 three-fold cross-validation. (d) Subsequently, we re-trained on the whole cohort and tested in an
228 external validation cohort. (e) Receiver operating characteristic (ROC, horizontal axis is 1-speci-
229 ficity and vertical axis is sensitivity) for HPV detection in TCGA-HNSC (N=412 patients), x marks
230 the performance of the current clinical gold standard (p16 immunohistochemistry) (f) ROC for
231 EBV detection in TCGA-STAD (N=317 patients), x marks the performance of EBV-encoded RNA
232 (EBER) in-situ hybridization, the diagnostic gold standard. Each ROC curve corresponds to one
233 cross-validation run.



234

235 **Figure 2: Tile-level classification yields high patient-level performance.** (a) Schematic of the pre-
 236 diction process: in a histological whole slide image, non-overlapping areas of 256 x 256 μm
 237 (“blocks” or “tiles”) were used to predict virus status and subsequently pooled on a patient-level
 238 by fraction of positive tiles. Image credit for icon <https://twemoji.twitter.com> (b) For a subset of
 239 tiles in the KCCH-STAD validation cohort, the predicted EBV score is plotted for true EBV-positive
 240 (left) and EBV-negative tiles (right). A small random symmetric x-y-offset was added to each point
 241 for better visibility. While most positive tiles attained a high EBV score, prediction of EBV-nega-
 242 tive tiles was noisy. More examples and more information about tile preprocessing is available at
 243 <http://doi.org/10.5281/zenodo.3247009>. (c) Despite noisy tile-level classification, patient-level
 244 prediction reached a high classification accuracy with an AUC of 0.81 in the independent KCCH-
 245 STAD validation set. (d) Representative tiles from the top and bottom quantile of EBV predictions.
 246 False positive tiles are lymphocyte-rich, which is a hallmark of virus-driven cancer and thus makes
 247 these misclassifications plausible.



248

249 **Figure 3: Feature visualization of viral morphological signatures in histological images by Deep**
250 **Dream.** (a) We used a modified VGG19 deep neural network that reads images in through the
251 input layer and outputs predictions in a two-neuron output layer. (b) Information flow from left
252 to right is used to classify images. The Deep dream algorithm uses the reverse direction to iteratively
253 create pseudo images for each output neuron. (c) Example of Deep dream pseudo-images
254 for HPV negative and positive HNSC with subjective manual description on morphological features
255 by expert pathologists, (b) corresponding images for EBV in STAD.

256

257 **Acknowledgements**

258 The results are in part based upon data generated by the TCGA Research Network: <http://cancerge->
259 [nome.nih.gov/](http://cancergenome.nih.gov/). Our funding sources are as follows. J.N.K.: RWTH University Aachen (START 2018-691906).
260 S. v. S.: RWTH University Aachen (START 2017-691742) and DFG (GRK 2375/1). P. Boor: DFG (SFB-
261 TRR57/P06, P25, M01, SFB-TRR219, LU 1360/3-1, BO 3755/3-1 and 6-1), BMBF (01GM1901A) and IZKF
262 (O3-7). T.L.: Horizon 2020 through the European Research Council (ERC) Consolidator Grant PhaseControl
263 (771083), Mildred-Scheel-Endowed Professorship from the German Cancer Aid, DFG (SFB-TRR57/P06, LU
264 1360/3-1), Ernst-Jung-Foundation Hamburg and IZKF (interdisciplinary center of clinical research) at
265 RWTH Aachen. A.T.P.: NIH/NIDCR (#K08-DE026500), Institutional Research Grant (#IRG-16-222-56) from
266 the American Cancer Society, and the University of Chicago Medicine Comprehensive Cancer Center Sup-
267 port Grant (#P30-CA14599). The authors would also like to thank the generous support from Chef Grant
268 Achatz, Chef Giuseppe Tentori, and Nick Kokonas.

269 **Competing interests**

270 The authors declare that no competing interests exist.

271

272 Bibliography

- 273 1. Bansal, A., Singh, M.P. & Rai, B. Human papillomavirus-associated cancers: A growing global
274 problem. *Int J Appl Basic Med Res* **6**, 84-89 (2016).
- 275 2. Mesri, E.A., Feitelson, M.A. & Munger, K. Human viral oncogenesis: a cancer hallmarks analysis.
276 *Cell Host Microbe* **15**, 266-282 (2014).
- 277 3. Chaturvedi, A.K., *et al.* Human papillomavirus and rising oropharyngeal cancer incidence in the
278 United States. *J Clin Oncol* **29**, 4294-4301 (2011).
- 279 4. Faraji, F., *et al.* The prevalence of human papillomavirus in oropharyngeal cancer is increasing
280 regardless of sex or race, and the influence of sex and race on survival is modified by human
281 papillomavirus tumor status. *Cancer* **125**, 761-769 (2019).
- 282 5. O'Sullivan, B., *et al.* Development and validation of a staging system for HPV-related
283 oropharyngeal cancer by the International Collaboration on Oropharyngeal cancer Network for
284 Staging (ICON-S): a multicentre cohort study. *Lancet Oncol* **17**, 440-451 (2016).
- 285 6. Kim, S.T., *et al.* Comprehensive molecular characterization of clinical responses to PD-1
286 inhibition in metastatic gastric cancer. *Nature Medicine* **24**, 1449-1458 (2018).
- 287 7. Jouhi, L., Hagstrom, J., Atula, T. & Makitie, A. Is p16 an adequate surrogate for human
288 papillomavirus status determination? *Curr Opin Otolaryngol Head Neck Surg* **25**, 108-112 (2017).
- 289 8. Coudray, N., *et al.* Classification and mutation prediction from non-small cell lung cancer
290 histopathology images using deep learning. *Nat Med* **24**, 1559-1567 (2018).
- 291 9. Kather, J.N., *et al.* Predicting survival from colorectal cancer histology slides using deep learning:
292 A retrospective multicenter study. *PLOS Medicine* **16**, e1002730 (2019).
- 293 10. Kather, J.N., *et al.* Deep learning can predict microsatellite instability directly from histology in
294 gastrointestinal cancer. *Nature Medicine* (2019).
- 295 11. The Cancer Genome Atlas, N., *et al.* Comprehensive genomic characterization of head and neck
296 squamous cell carcinomas. *Nature* **517**, 576 (2015).
- 297 12. The Cancer Genome Atlas Network, *et al.* Comprehensive molecular characterization of gastric
298 adenocarcinoma. *Nature* **513**, 202 (2014).
- 299 13. Aoyama, T., *et al.* Identification of a high-risk subtype of intestinal-type Japanese gastric cancer
300 by quantitative measurement of the luminal tumor proportion. *Cancer Med* **7**, 4914-4923
301 (2018).
- 302 14. Campbell, J.D., *et al.* Genomic, Pathway Network, and Immunologic Features Distinguishing
303 Squamous Carcinomas. *Cell Rep* **23**, 194-212 e196 (2018).
- 304 15. Kostic, A.D., *et al.* PathSeq: software to identify or discover microbes by deep sequencing of
305 human tissue. **29**, 393 (2011).
- 306 16. The Cancer Genome Atlas Research, N., *et al.* Integrated genomic and molecular
307 characterization of cervical cancer. *Nature* **543**, 378 (2017).
- 308 17. Liu, Y., *et al.* Comparative Molecular Analysis of Gastrointestinal Adenocarcinomas. *Cancer Cell*
309 **33**, 721-735 e728 (2018).
- 310 18. Hewitt, L.C., *et al.* Epstein-Barr virus and mismatch repair deficiency status differ between
311 oesophageal and gastric cancer: A large multi-centre study. *Eur J Cancer* **94**, 104-114 (2018).
- 312 19. Bankhead, P., *et al.* QuPath: Open source software for digital pathology image analysis. *Scientific*
313 *reports* **7**, 16878 (2017).
- 314 20. Simonyan, K. & Zisserman, A. Very deep convolutional networks for large-scale image
315 recognition. (2014).
- 316 21. Mori, M., *et al.* Epstein-Barr virus-associated carcinomas of the esophagus and stomach.
317 *Archives of pathology & laboratory medicine* **118**, 998-1001 (1994).

- 318 22. Westra, W.H. The Morphologic Profile of HPV-Related Head and Neck Squamous Carcinoma:
319 Implications for Diagnosis, Prognosis, and Clinical Management. *Head and Neck Pathology* **6**, 48-
320 54 (2012).
- 321 23. Mordvintsev, A., Olah, C. & Tyka, M. Inceptionism: Going Deeper into Neural Networks. in
322 *Google AI Blog* (2015).
- 323 24. Kochanny, S., *et al.* High-accuracy HPV testing versus p16 IHC using multiple clinically relevant
324 outcomes: The University of Chicago Experience. *Journal of Clinical Oncology* **36**, 6020-6020
325 (2018).
- 326 25. Camargo, M.C., *et al.* Validation and calibration of next-generation sequencing to identify
327 Epstein-Barr virus-positive gastric cancer in The Cancer Genome Atlas. *Gastric Cancer* **19**, 676-
328 681 (2016).
- 329 26. Schaumberg, A.J., Rubin, M.A. & Fuchs, T.J. H&E-stained Whole Slide Image Deep Learning
330 Predicts SPOP Mutation State in Prostate Cancer. 064279 (2018).
- 331 27. Carter, S., Armstrong, Z., Schubert, L., Johnson, I. & Olah, C. Activation Atlas. *Distill* (2019).

332