

PanNuke Dataset Extension, Insights and Baselines

Jevgenij Gamper*, Navid Alemi Koohbanani*, Ksenija Benes, Simon Graham, Mostafa Jahanifar, Seyyed Ali Khurram, Ayesha Azam, Katherine Hewitt and Nasir Rajpoot

Abstract—The emerging area of computational pathology (CPath) is ripe ground for the application of deep learning (DL) methods to healthcare due to the sheer volume of raw pixel data in whole-slide images (WSIs) of cancerous tissue slides. However, it is imperative for the DL algorithms relying on nuclei-level details to be able to cope with data from ‘the clinical wild’, which tends to be quite challenging. We study, and extend recently released PanNuke dataset consisting of nearly 200,000 nuclei categorized into 5 clinically important classes for the challenging tasks of segmenting and classifying nuclei in WSIs. Previous pan-cancer datasets consisted of only up to 9 different tissues and up to 21,000 unlabeled nuclei [1] and just over 24,000 labeled nuclei with segmentation masks [2]. PanNuke consists of 19 different tissue types that have been semi-automatically annotated and quality controlled by clinical pathologists, leading to a dataset with statistics similar to ‘the clinical wild’ and with minimal selection bias. We study the performance of segmentation and classification models when applied to the proposed dataset and demonstrate the application of models trained on PanNuke to whole-slide images. We provide comprehensive statistics about the dataset and outline recommendations and research directions to address the limitations of existing DL tools when applied to real-world CPath applications.

I. INTRODUCTION

The success of convolutional neural networks (CNNs) in computer vision (CV) algorithms applied to natural image and medical imaging tasks can generally be attributed to the availability of large datasets and computing power [3]–[7]. Given the excellent performance on these tasks, measured by an *average* metric evaluated over a large dataset, CNNs have sparked hope and promise in healthcare applications [8], [9].

The field of computational pathology (CPath) is witnessing a rapid rise in the research and development of deep learning (DL) models for quantitative profiling of spatial patterns in digitized whole-slide images (WSIs) of cancerous tissue slides that are rich in content and information [10]. Numerous studies have demonstrated the potential of deep learning (DL) models in detecting cancer, classifying tissue, identifying diagnostically relevant structures and even inferring genetic sub-types [9], [11]–[15].

* First authors contributed equally.

J.Gamper, N.A.Koohbanani, S.Graham and N.Rajpoot are with the Department of Computer Science, University of Warwick, UK.

J.Gamper and S.Graham are also with the Mathematics for Real-World Systems Centre for Doctoral Training, University of Warwick, UK.

M.Jahanifar is with the Department of Research and Development, NRP Co., Tehran, Iran

K.Benes is with the Department of Pathology, Royal Wolverhampton NHS Trust, UK

S.A.Khurram is with the Department of Clinical Dentistry, University of Sheffield, UK

A.Azam and K.Hewitt are with the Department of Pathology at University Hospitals Coventry and Warwickshire, Coventry, UK

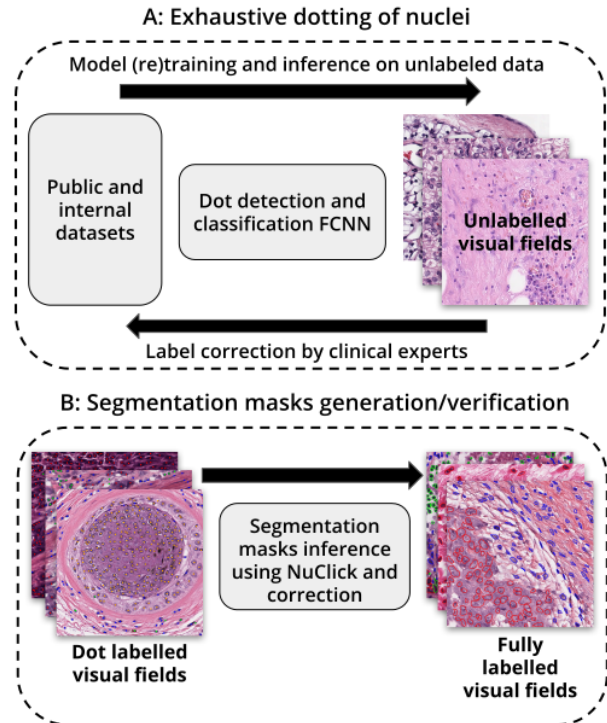


Fig. 1: Illustration of PanNuke label generation and verification.

It would be fair to state that challenge contests have become a popular mean for attracting attention to a particular dataset or a task in medical imaging. In computational pathology, for instance, a variety of deep CNNs for automated nucleus segmentation have been developed on a dataset consisting of 21,623 nuclei in a total of 32 images of the size $1,000 \times 1,000$ pixels and released as part of the MoNuSeg challenge contest [1]. However, the use and validity of results in most challenge contests is questionable due to the limited diversity [16]. PanNuke, released under the CC license is however a diverse pan-cancer dataset that has undergone clinical quality control (QC).

Second, CNN models applied to medical images are more powerful than the datasets that we apply them to, these models are known to suffer from over-fitting to surface statistical regularities [17]–[19]. A model can obtain stronger inductive biases (a set of assumptions that a model uses to predict on a certain task) through multi-task learning, and therefore be more robust in practice [20]. The current approach in computational pathology is to train a model for a specific tissue, or even just a specific disease sub-type classification. On the other hand, we argue for a top-down algorithm development in

medical imaging i.e. develop general and user-friendly tools that can easily be used by domain experts, and for bottom-up approach for dataset creation for supervised learning. Far too much medical imaging community’s effort is focused on semi-supervised and unsupervised studies, however, these are far from clinically applicable and would always under-perform compared to supervised approaches trained on accurately labeled ground truth data. Models trained on PanNuke, as we demonstrate in Figure 3, could be used to assist detailed semi-automatic labeling in 19 different tissues. The results of nuclei detection and classification can be used for tissue pheno-typing as demonstrated in Colon tissue by Javed *et al.* [12].

This work is motivated by the fact that publicly available nucleus segmentation and classification datasets do not often match the distribution of data in ‘the clinical wild’, as can be seen in Figure 2 which shows the results of a nucleus detection model [21] trained on the MoNuSeg challenge dataset [1]. It can be observed that when the model is applied to a few images that contain commonly found artifacts in clinical practice, there are several false detections which could lead to incorrect or misleading results in downstream CPath analysis. A similar phenomenon has been demonstrated by Oakden-Rayner *et al.* [22] in the literature and in a validation study of DL models applied to radiology images where model performance dropped significantly in a real-world environment [23].

The main contributions of this work are summarized as below:

- We present PanNuke¹, the largest and the most diverse to date dataset for nucleus segmentation and classification, that has been annotated in a semi-automated manner and quality-controlled by clinical professionals;
- We accelerate the process of verification (i.e., quality control) by the clinical professionals by incorporating NuClick [24] during the generation of segmentation mask, as shown in Figure 1;
- We evaluate the performance of several nucleus segmentation models on PanNuke, which is 26 times larger in terms of unique 224×224 training patches than the previous MoNuSeg challenge dataset [1];
- We provide a full schema that could be used by the algorithm developers and that applies to other tissues not included in the PanNuke dataset. We show that segmentation models trained on PanNuke generalize to tissues like brain that were not part of the dataset.
- By releasing PanNuke to the broader CV research community, we encourage the community to develop new DL models to help push forward clinically relevant research.

These are not the only goals of this research. In fact, we hope to draw attention to the established tendencies in the field and their systemic impact on the risks as well as positive outcomes of CV in the healthcare domain and its progress.

In the following sections, we describe the relevant literature, our methodology, the quality of the automatically generated nucleus segmentation masks, provide qualitative and quantita-

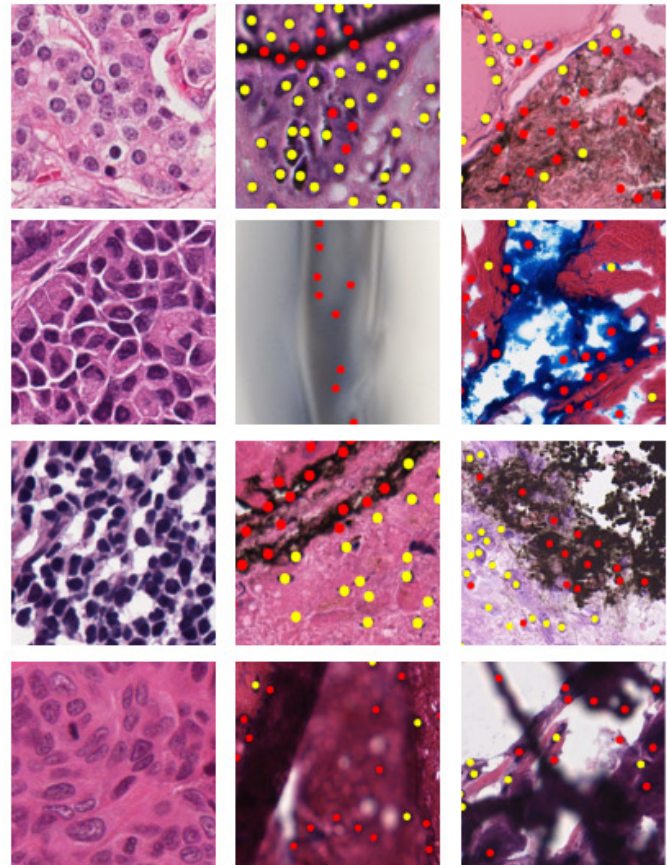


Fig. 2: 1st column: A selection of visual fields from the Kumar dataset *et al.* [25]. 2nd-3rd columns: A selection of visual fields in PanNuke with output of a detector trained on [25] overlaid on the images. False positives (shown as red dots, as opposed to true positives as yellow dots) are clearly visible in the areas of burnt tissue, blur or other tissue processing or scanning artifacts.

tive analysis of the dataset and its significance. Finally, we evaluate the performance of existing approaches to nucleus segmentation and classification.

A. Related Work

Recent work in computational pathology has demonstrated that nuclear features can be effectively used for: cancer scoring, bio-marker discovery, cancer recurrence prediction as well as for predicting treatment effectiveness [26]–[29]. However, these are small sample size studies limited to a single tissue, and are focused on commonly studied tissues such as lung, breast, prostate and colon, due to the lack of data for other tissues. Javed *et al.* [12] demonstrated that inferred nuclear categories may help in classifying different tissue phenotypes within colon slides by constructing a graph of nuclei connections and detecting communities. By training models on PanNuke and applying them to WSIs at scale, these studies could be extended to 19 different tissues.

Closest in scale and real-world proximity to PanNuke is the work of Hosseini *et al.* [30], who have provided multi-category patch-level annotations. Due to having only patch annotations,

¹Download dataset here: <https://warwick.ac.uk/fac/sci/dcs/research/tia/data/pannuke>

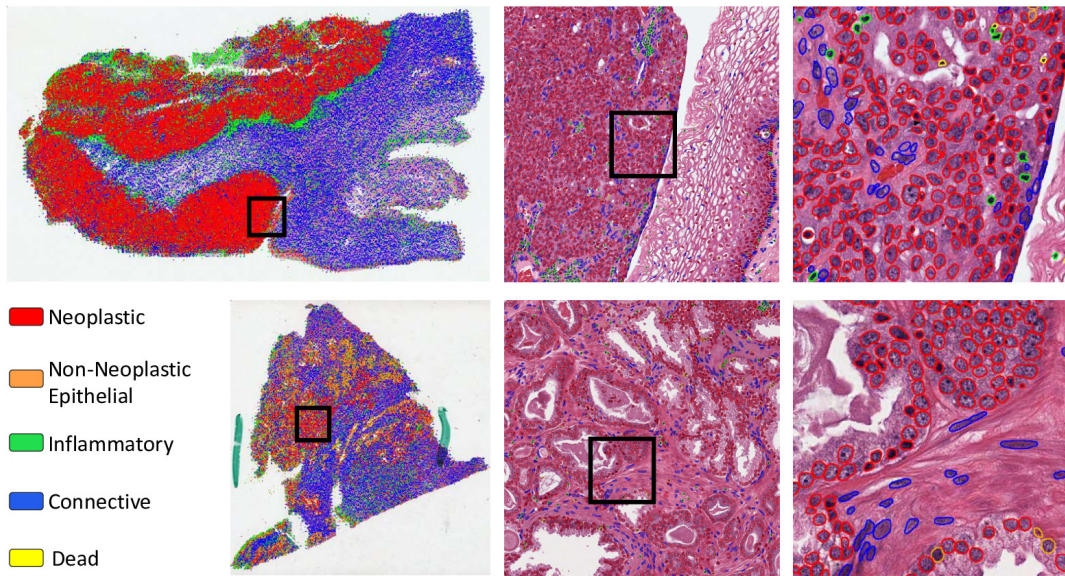


Fig. 3: Visualisation of applying nuclei segmentation and classification network trained on PanNuke to unseen whole-slide images. Top row: Cervix tissue, with a visible differentiation between tumor and other tissue types. Bottom row: Prostate tissue.

Chan *et al.* [31] pursued semantic segmentation using class activation maps that significantly under-performed pixel-wise supervised models on the GlaS challenge dataset [32]. Across the board, large CPath studies have been somewhat limited by the lack of granular annotations [3], [33], [34]. However, PanNuke provides pixel-level boundary annotation for every individual nucleus, a building block of any organ’s tissue. As a result, using models trained on PanNuke further semi-automatic labeling of tumor or tissue phenotypes is feasible [12] (Figure 3).

II. THE PANNUKE DATASET

In the sections below, we describe our methodology, discuss the quality of the automatically generated nucleus segmentation masks and provide qualitative and quantitative analysis of the dataset and its significance in algorithmic and practical terms.

A. Dataset Generation

Data Labeling: First, we aggregated a set of publicly available nucleus classification and detection datasets to create an initial dataset for semi-automatic ground truth generation. For this, we trained a fully convolutional neural network (FCNN) for nucleus detection using 4 publicly available datasets: Kumar [25], CPM2017 [36], 15 visual field from TCGA [37] that we have labeled ourselves, and a dataset of bone marrow visual fields [38]. Kumar is a dataset of 16 visual

fields consisting of 7 different tissue types (liver, prostate, kidney, breast, stomach, colorectal and bladder) and CPM2017 is a dataset of 64 visual fields of 4 cancer types (glioblastoma, low grade glioma, head neck squamous cell carcinoma and non-small cell lung cancer). Therefore, in total we initially utilize 106 visual fields. We then trained a convolutional neural network on nuclei CNN patches (see Figure 4 for illustration) extracted from the datasets described in Table 1², that were re-labeled according to the PanNuke categories. The schema for nuclei categories will be described in the subsequent sections. Using the CNN, we exhaustively classified all of the detected nuclei in the above mentioned datasets, which were then verified by a team of expert pathologists. After, we sampled 2,000 visual fields from more than 20,000 WSIs in 19 different tissues obtained from TCGA and also from a local hospital. Random sampling of visual fields allowed us to address the selection bias present in available public datasets, in fact visual fields with common clinical artifacts demonstrated in Figure 2 are from the final PanNuke dataset.

When sampling visual fields of tissue from WSIs, the tissue may have been originally frozen or paraffin embedded and also the WSIs may have been scanned with a maximum resolution of either 20 \times or 40 \times . We re-sized the selected visual fields so that all images were at 40 \times resolution and we excluded frozen tissue from the study. We then re-sampled the visual fields so that the images present in the dataset were diagnostically

²BrestPathQ source: <https://breastpathq.grand-challenge.org>

TABLE I: Data used to initialize semi-automatic labeling.

	PanNuke Nuclei Categories						Total
	Neoplastic	Non-Neo Epithelial	Inflammatory	Connective	Dead	Non-Nuclei	
MoNuSeg	5,927	836	1,698	906	0	0	9,367
Colon Nuclei	4,685	7,544	6,003	4,468	2,547	0	25,247
BreastPathQ	9,802	0	2,139	0	0	0	11,941
Nuclei Attribute	0	0	0	0	0	500	500
Total	20,414	8,380	9,840	5,374	2,547	500	47,055

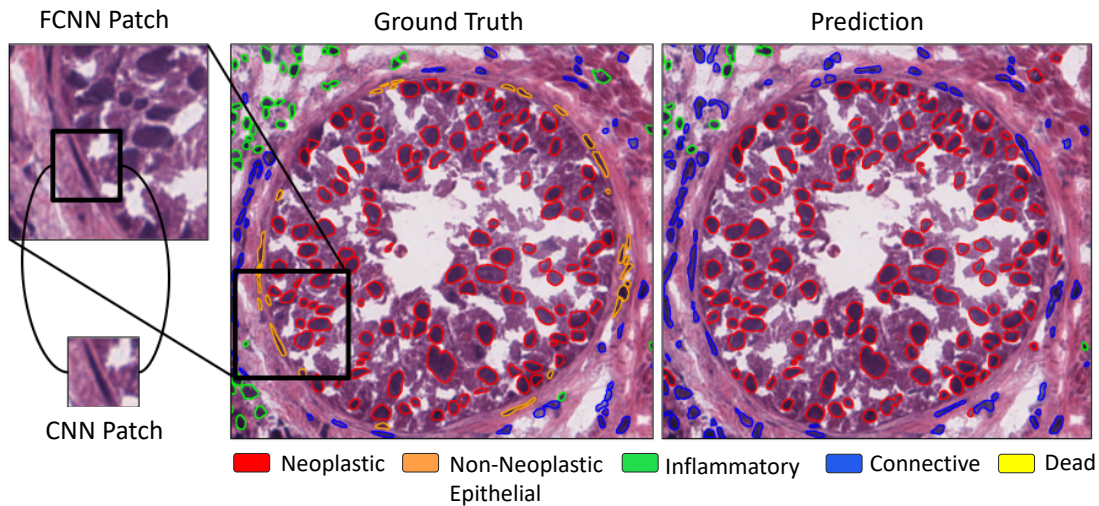


Fig. 4: Ground truth labels and segmentation masks verified by pathologists alongside model prediction for bladder tissue visual field. FCNN patch represents a 224×224 patch commonly used for training fully convolution segmentation models, right below it is a patch used by [35] for training a CNN to classify each individual nucleus.

relevant and also reflected the true variation of tissue in each organ. We also ensured that artifacts remained in the dataset because they are inevitable when facing data in the ‘clinical wild’.

We then proceeded to Part A of annotation process, as illustrated in Figure 1. Iterating 7 times, where after each stage pathologists would verify and re-label the detected and classified nuclei dots. In each stage, the FCNN detection and classification model was trained with the new collected annotations to provide better predictions for next iteration. Eventually, this leads to a dataset of 481 visual fields with a total of 189,744 exhaustively annotated nuclei, verified by domain experts.

Mask Generation: For the final version of PanNuke we used NuClick [24]. NuClick is a method that enables accurate segmentation mask generation from a single point. Therefore, this enables us to produce many segmentation masks at a relatively low cost and also eases the verification process because only a single point is required from the pathologist, as opposed to the entire mask. Figure 5 presents a selection of masks generated by a segmentation FCNN, as well as by NuClick. In the first row, the FCNN has mistaken pigment in skin tissue as nuclei. However, because NuClick is conditioned on verified nuclei dots, only nuclei pixels are segmented. Also, the FCNN frequently misses elongated nuclei compared to NuClick and often struggles to segment nuclei with indistinct boundaries, as pictured in the second column and third row. Finally, the last row shows that FCNN only segments the

nucleolus, whereas NuClick segments the entire nucleus.

Using the proposed semi-automatic pipeline, we generated and verified 189,744 nuclei from more than 20,000 WSIs. Some examples of generated ground truth are depicted in the Figure 6.

B. Dataset Description and Statistics

Dataset Schema: For the purpose of this work, we derived a *schema* in Table II for generating nuclei labels that is clinically sound and is shared across the 19 tissues within the dataset. The proposed schema provides insight into how we categorized all nuclei and can be used to appropriately subtype nuclei in future studies. Our schema is consistent with the nuclei categories used in previous tissue-specific studies [2], [30], [35] and therefore work previously developed on these datasets can be seamlessly extended to PanNuke. We split the cell types between Neoplastic and Non-neoplastic cells.

Neoplasm (‘new growth’) includes any tumor, malignant or benign. It includes carcinomas, sarcomas, melanomas, lymphomas, etc. These are all tumors but originate from different cell types: carcinomas from epithelial; sarcomas from soft tissue; melanomas from melanocytes; lymphomas from lymphoid cells and so on. As such, all tumorous cells in PanNuke are labeled as Neoplastic.

Non-neoplastic covers everything else, from normal to inflammatory, degenerative, metaplastic, atypia etc. For the purpose of this exercise, atypia is under the heading of non-neoplastic, although some researchers may argue that they

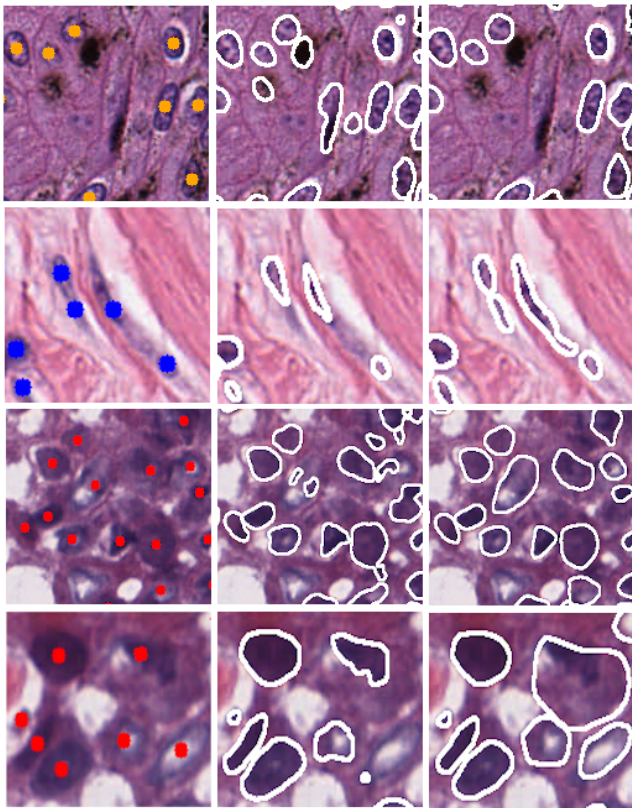


Fig. 5: Left column: Pathologist verified nuclei dots; Middle: CNN generated segmentation masks; Right: NuClick generated segmentation masks and subsequently verified. Color dots are consistent with the legend in Figure 4.

have clonal changes and potential for turning into neoplastic cells. As such, non-neoplastic labels in PanNuke are: epithelial; connective/soft tissue cells; inflammatory and dead cells. Here, connective tissue cells have the potential to become neoplastic, whereas inflammatory cells typically cannot become neoplastic. Inflammatory cells include lymphoid and macrophage cells in PanNuke. Dead cells can arise from either neoplastic or non-neoplastic cells, but in this study we refer to them as non-neoplastic.

Dataset Statistics: In general, most nuclei types that we provide in our schema are represented in all tissues considered in PanNuke, but the distribution of the nuclei count per class may vary from tissue to tissue. This can be seen in Figure 7, where we observe that the total nuclei count per tissue as well as per class changes between tissue types.

In our experience, a pathologist when annotating nuclei frequently refers to the WSI at a lower resolution, to observe the surrounding structures. Seminal work in automatic nuclei classification by Sirinukunwattana *et al. et al.* [35] considered only nuclear patches when performing classification (CNN patch in Figure 4), or small image patches containing a few nuclei in the more recent work using fully convolutional networks (FCNN patch in Figure 4) - both approaches have demonstrated high accuracy on average. What does it mean for us? Either that pathologist overfits given the surrounding; or that looking at a nuclear patch to classify nuclei is not

TABLE II: Nucleus classification schema for computational pathology.

#	Level-1	Level-2	Level-3
1	Epithelial	Neoplastic	Malignant ⁺
2			Benign ⁻
3		Non neoplastic	
4	Connective/ Soft tissue cells		Fibroblasts
5			Endothelial
6			Myo-fibroblasts
7			Fibers
8			Adipocytes
9			Trombocytes
10	Lympho-reticular cells	Leukocytes	Erythrocytes
11			Eosinophils
12			Basophils
13		Lymphoid cells	Neutrophils
14			Plasma cells
15			Lymphocytes
16			Macrophages/Histiocytes
17			Mast cells
18	Nervous system cells	CNS	Oligo-dendrocytes
19			Microglia
20			Astrocyte
21			Ependymal cells
22		Peripheral Nervous	Schwann cells
23		Ganglia	Ganglion cells
24	Dead		Apoptotic
25			Necrotic

sufficient; or an *average* accuracy metric is not a clinically relevant measure of algorithmic performance. Recently Oaken *et al.* [39] has demonstrated effects of *hidden stratification* in medical imaging, i.e. when labels used in CV or ML studies do not represent clinical reality, and that *average* performance is not a strong measure of applicability. As Figure 4 demonstrates *hidden stratification* is also present in histology, if one is to use nuclei classification for downstream tasks. Carcinoma (malignant epithelial tumor/neoplasm) can be non invasive (i.e. in situ) or invasive. They are both composed of malignant cells but in situ is still "bounded" by either basal cell layer/myoepithelial layer or basement membrane (depending on organ site) hence called non-invasive. Invasive carcinomas have lost basal cells/myoepithelial cells. So in theory, non invasive cancer should have no or low potential for lympho-vascular space invasion, distant metastases and similar (i.e. they should have better prognosis than or may not need as radical treatment as invasive carcinomas). For this specific image, for example, it would appear that this tumor cells are 'bounded' by basement membrane cells (Figure 4 Ground Truth). While FCNN based model (Figure 4 Prediction) classifies epithelial cells surrounding tumor cells as part of connective tissue category, evidently due to their shapes and the tendency of deep models to look for simplest association that describe the relationship between the feature and target on average [17], [19]. Notably, these cells would be unlikely correctly classified by ANY algorithm. We either need to incorporate prior knowledge about the underlying mechanisms in the tissue in order to solve these tasks as precisely as a pathologist, or acquire multiple, similar cases and train on substantially larger image sizes than FCNN patch.

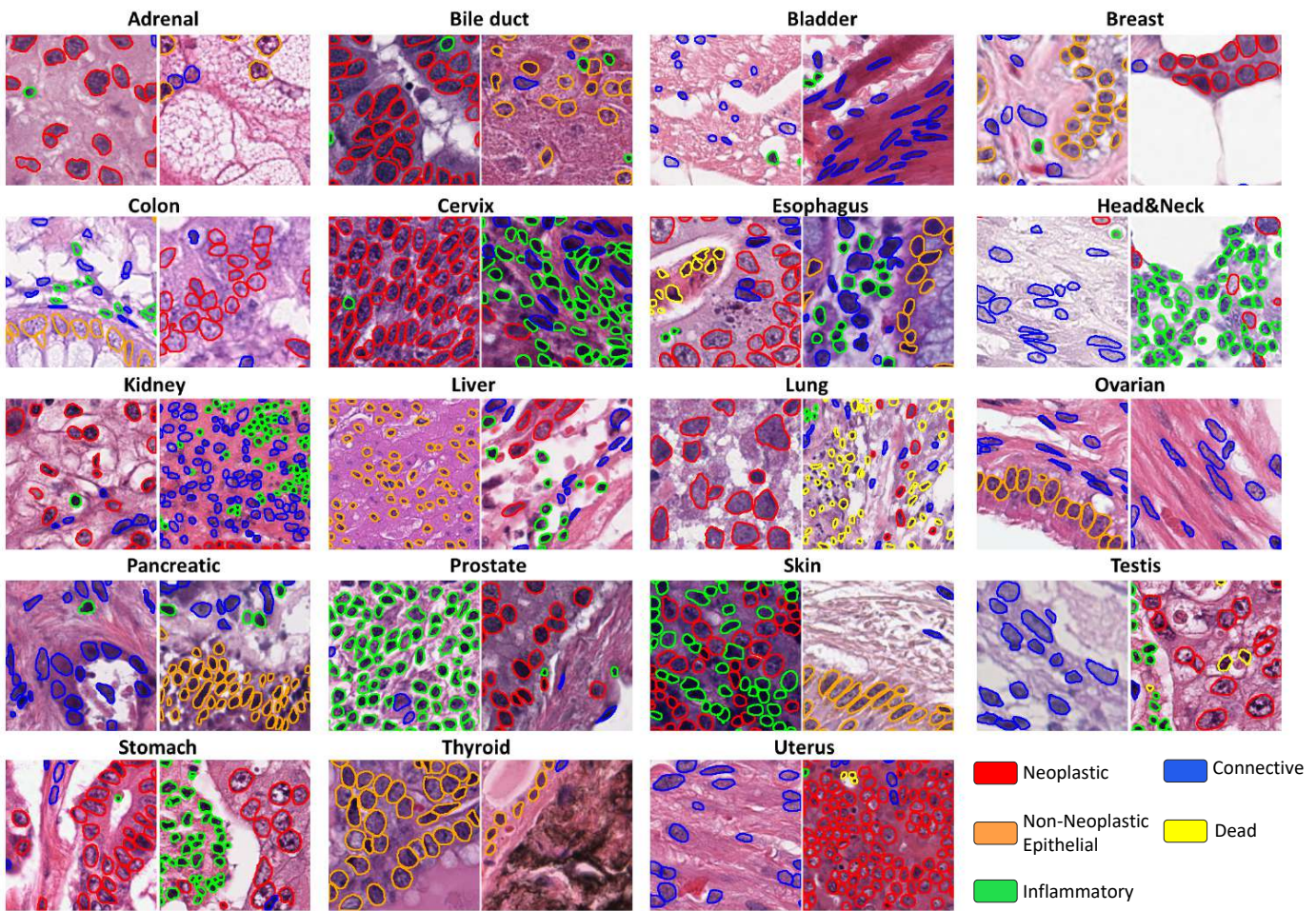


Fig. 6: Example of PanNuke patches and their ground truth annotations overlaid.

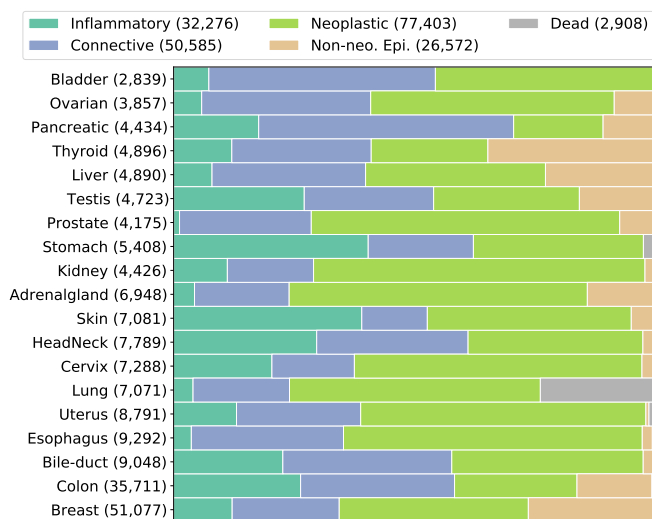


Fig. 7: A comparative plot of class distributions per tissue. Numbers in parenthesis represent the total number of nuclei within that category or tissue type.

These are just a few examples of cases which are difficult to simply classify/diagnose as neoplastic or non-neoplastic. Other challenging examples could not be included given the limited space, but in clinical practice these cases are more frequent as compared to how they are portrayed in the AI/digital pathology literature.

III. PERFORMANCE BENCHMARKS

Previous nucleus segmentation datasets provided visual fields and therefore patch extraction between different methods for training and testing is not standardized. For PanNuke we pre-extract patches and split into 3 randomized training, validation and testing folds for a fair model comparison. For every fold we split every tissue into three sections by ensuring that each contains an equal portion of the smallest class within it (refer to Figure 7). Here, we apply recent and well known models on PanNuke to create a benchmark for further research using this dataset.

A. Baseline Models

There are not an abundance of models that perform simultaneous segmentation and classification and therefore we adapt several top-performing instance segmentation models so that they additionally classify each nucleus. We quantified

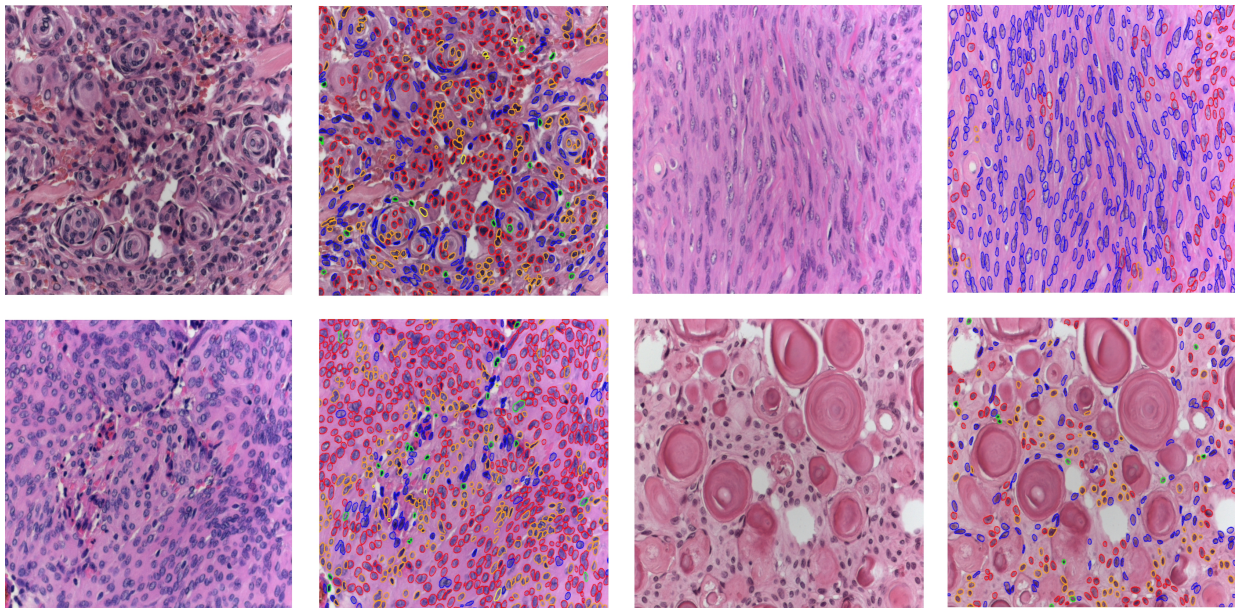


Fig. 8: Brain tissue visual fields from Qureshi *et al.* [40] and predictions overlay using HoVer-Net trained on PanNuke.

the performance of 4 models on PanNuke: DIST [41] which utilizes the distance map of instances as the target and we add another branch for semantic pixel-wise classification; MaskRCNN [42] which is a state-of-the-art instance segmentation network for natural images; Micro-Net [43] which was proposed for nuclear and gland segmentation and HoVer-Net [2], which uses the concept of horizontal and vertical distance maps to separate clustered nuclei. HoVer-Net does not need to be adapted as it inherently performs simultaneous nuclear instance segmentation and classification. Note, that each of the above mentioned models performed extensive comparison with competing segmentation methods and therefore we focus on the best performing models.

In addition to the segmentation models, we also implemented a detection-based U-Net, similar to the model proposed by Xie *et al.* [44], where the detection maps for each class were used as the target for training the network. Specifically, the detection map for each nucleus was converted to a 2D Gaussian centered at the true nucleus centroid.

B. Evaluation

Instance Segmentation: To quantify the instance segmentation performance of each of the models trained on PanNuke, we use panoptic quality (PQ) [2], [45]. An in depth discussion as to why we choose to utilize PQ over other recently used metrics for nucleus segmentation is discussed by Graham *et al.* [2]. Specifically, we used multi-class PQ (mPQ) and binary PQ (bPQ) that assumes that all nuclei belong to one class. For mPQ, the PQ is calculated independently for each positive class and then the results are averaged. Therefore, the metric is insensitive to class imbalance. Our main criterion for evaluating model performance is the average mPQ over all of the tissues, which therefore equally weights the contribution of each tissue type. The IoU threshold for determining a true positive during PQ calculation is set to 0.5. As part of

this work, we provide the implementation of our evaluation framework³ to encourage reproducibility, however this is a preliminary release of the code and we are working on a full repository release.

To provide insight into how each model performs for different types of nuclei, in Table III, we report mPQ and bPQ for all 19 tissue types separately. We observe that HoVer-Net achieves the best performance for most tissue types and is reflected by the greatest average score over all tissues for mPQ and bPQ. We also report PQ for each type of nucleus in Table IV. Dead cells obtain a low PQ for all models because these nuclei are very small and therefore achieving an IoU>0.5 (PQ criterion for true positive) is difficult. In some cases, distinguishing between neoplastic and non-neoplastic nuclei proved to be challenging, yet it must be emphasized that this also can be a challenging task for the pathologist, where they often need to assess contextual information before confirming the the type of nucleus. As deduced from Figure 7 class imbalance may also lead to poor performance for dead cell and non-neoplastic classes.

Detection: In order to allow cross comparison with detection models (as opposed to segmentation), we reported the F_1 , precision and recall for the overall detection quality in Table V. Here, a true positive was considered as a detection within 12 pixels of the labeled centroid [35]. We also calculated. In order to report the detection performance for segmentation models, we extracted the centroids of each instance as detection points. We observed that segmentation models generally performed better than the detection model. We hypothesize that this is because the detection-based model does not incorporate boundary information.

³Evaluation code: <https://github.com/TIA-Lab/PanNuke-metrics>

TABLE III: Average mPQ and bPQ across three dataset splits. We also provide the standard deviation (STD) across these splits in the final row.

	DIST		Mask-RCNN		Micro-Net		HoVer-Net	
	mPQ	bPQ	mPQ	bPQ	mPQ	bPQ	mPQ	bPQ
Adrenal Gland	0.3442	0.5603	0.3470	0.5546	0.4153	0.6440	0.4812	0.6962
Bile Duct	0.3614	0.5384	0.3536	0.5567	0.4124	0.6232	0.4714	0.6696
Bladder	0.4463	0.5625	0.5065	0.6049	0.5357	0.6488	0.5792	0.7031
Breast	0.3790	0.5466	0.3882	0.5574	0.4407	0.6029	0.4902	0.6470
Cervix	0.3371	0.5309	0.3402	0.5483	0.3795	0.6101	0.4438	0.6652
Colon	0.2989	0.4508	0.3122	0.4603	0.3414	0.4972	0.4095	0.5575
Esophagus	0.3942	0.5295	0.4311	0.5691	0.4668	0.6011	0.5085	0.6427
Head & Neck	0.3177	0.4764	0.3946	0.5457	0.3668	0.5242	0.4530	0.6331
Kidney	0.3339	0.5727	0.3553	0.5092	0.4165	0.6321	0.4424	0.6836
Liver	0.3441	0.5818	0.4103	0.6085	0.4365	0.6666	0.4974	0.7248
Lung	0.2809	0.4978	0.3182	0.5134	0.3370	0.5588	0.4004	0.6302
Ovarian	0.3789	0.5289	0.4337	0.5784	0.4387	0.6013	0.4863	0.6309
Pancreatic	0.3395	0.5343	0.3624	0.5460	0.4041	0.6074	0.4600	0.6491
Prostate	0.3810	0.5442	0.3959	0.5789	0.4341	0.6049	0.5101	0.6615
Skin	0.2627	0.5080	0.2665	0.5021	0.3223	0.5817	0.3429	0.6234
Stomach	0.3369	0.5553	0.3684	0.5976	0.3872	0.6293	0.4726	0.6886
Testis	0.3278	0.5548	0.3512	0.5420	0.4088	0.6300	0.4754	0.6890
Thyroid	0.2574	0.5596	0.3037	0.5712	0.3712	0.6555	0.4315	0.6983
Uterus	0.3487	0.5246	0.3683	0.5589	0.3965	0.5821	0.4393	0.6393
Average across tissues	0.3406	0.5346	0.3688	0.5528	0.4059	0.6053	0.4629	0.6596
STD across splits	0.0156	0.00975	0.00465	0.00762	0.00816	0.00499	0.00758	0.00364

TABLE IV: Average PQ across three dataset splits for each nuclear category.

	Neo	Non-Neo Epi	Inflam	Conn	Dead
DIST	0.439	0.290	0.343	0.275	0.000
Mask-RCNN	0.472	0.403	0.290	0.300	0.069
Micro-Net	0.504	0.442	0.333	0.334	0.051
HoVer-Net	0.551	0.491	0.417	0.388	0.139

C. Generalisation to other tissues

We speculated that models trained on PanNuke would likely generalise to other tissues and applied the best performing model to brain tissue as demonstrated in Figure 8. Within this tissue, the model was able to perform a successful segmentation of all nuclei, but found it challenging to predict the correct nuclear categories. The algorithm trained with PanNuke performs favourably for segmentation of nuclei for the 4 images in Figure S3 (DICE value 0.796, mPQ 0.28 and bPQ 0.51) from a completely unseen source (Germany) and tissue type (brain), as compared to the tissue-wise average in Table 3.

IV. CONCLUDING REMARKS

In this paper, we presented a semi-annotated and quality-controlled dataset with detailed boundaries and class labels for 5 main types of nuclei for multiple different cancerous tissue types. This work is motivated by the observation that the use and validity of results in most challenge contests is questionable due to the limited nature of challenge datasets [16]. For example, even on ImageNet [46] (which happens to be many orders or magnitude larger than [25]), there is evidence of overfitting due to multiple-hypothesis testing in architecture development [47]. In addition to selection bias in the available datasets, results are reported on labels that do not always meaningfully describe the variation within the population. This work, while providing a significant contribution in

modeling and dataset size compared to any previous work, is only a small step in the direction of safe and robust application of CV in CPath. Similarly to Esteva *et al.* [48], we offer a careful treatment of PanNuke labels, discuss the real world complexities of the task and offer *schemas* that researchers can use to push nuclei classification research further.

REFERENCES

- [1] N. Kumar, R. Verma, D. Anand, Y. Zhou, O. F. Onder, E. Tsougenis, H. Chen, P. A. Heng, J. Li, Z. Hu, *et al.*, "A multi-organ nucleus segmentation challenge," *IEEE transactions on medical imaging*, 2019.
- [2] S. Graham, Q. D. Vu, S. E. A. Raza, A. Azam, Y. W. Tsang, J. T. Kwak, and N. Rajpoot, "Hover-net: Simultaneous segmentation and classification of nuclei in multi-tissue histology images," *Medical Image Analysis*, vol. 58, p. 101563, 2019.
- [3] G. Campanella, V. W. K. Silva, and T. J. Fuchs, "Terabyte-scale Deep Multiple Instance Learning for Classification and Localization in Pathology," *arXiv:1805.06983 [cs]*, May 2018. arXiv: 1805.06983.
- [4] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet Classification with Deep Convolutional Neural Networks," in *Advances in Neural Information Processing Systems 25* (F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, eds.), pp. 1097–1105. Curran Associates, Inc., 2012.
- [5] Y. Liu, K. Gadepalli, M. Norouzi, G. E. Dahl, T. Kohlberger, A. Boyko, S. Venugopalan, A. Timofeev, P. Q. Nelson, G. S. Corrado, J. D. Hipp, L. Peng, and M. C. Stumpe, "Detecting Cancer Metastases on Gigapixel Pathology Images," *arXiv:1703.02442 [cs]*, Mar. 2017. arXiv: 1703.02442.
- [6] K. Sirinukunwattana, J. P. W. Pluim, H. Chen, X. Qi, P.-A. Heng, Y. B. Guo, L. Y. Wang, B. J. Matuszewski, E. Bruni, U. Sanchez, A. Bhm, O. Ronneberger, B. B. Cheikh, D. Racoceanu, P. Kainz, M. Pfeiffer, M. Urschler, D. R. J. Snead, and N. M. Rajpoot, "Gland Segmentation in Colon Histology Images: The GlaS Challenge Contest," *arXiv:1603.00275 [cs]*, Mar. 2016. arXiv: 1603.00275.
- [7] A. J. Schaumberg, W. Juarez, S. J. Choudhury, L. G. Pastrian, B. S. Pritt, M. Prieto Pozuelo, R. Sotillo Sanchez, K. Ho, N. Zahra, B. D. Sener, S. Yip, B. Xu, S. R. Annavarapu, A. Morini, K. A. Jones, K. Rosado-Orozco, S. J. Sirintrapun, M. Aly, and T. J. Fuchs, "Large-Scale Annotation of Histopathology Images from Social Media," Aug. 2018.
- [8] A. Esteva, A. Robicquet, B. Ramsundar, V. Kuleshov, M. DePristo, K. Chou, C. Cui, G. Corrado, S. Thrun, and J. Dean, "A guide to deep learning in healthcare," *Nature Medicine*, vol. 25, pp. 24–29, Jan. 2019.

TABLE V: Precision (P), Recall (R) and F_1 score for detection and classification. All results are the average across three dataset splits. For segmentation predictions, the centroid of each nucleus is extracted for computing these metrics.

Detection				Classification														
				Neoplastic			Non-Neo Epithelial			Inflammatory			Connective			Dead		
	P	R	Fd	P	R	F1	P	R	F1	P	R	F1	P	R	F1	P	R	F1
Det U-Net	0.73	0.59	0.65	0.40	0.47	0.43	0.27	0.31	0.29	0.32	0.45	0.37	0.34	0.38	0.36	0.00	0.00	0.00
DIST	0.74	0.71	0.73	0.49	0.55	0.50	0.38	0.33	0.35	0.42	0.45	0.42	0.42	0.37	0.39	0.00	0.00	0.00
MRCNN	0.76	0.68	0.72	0.55	0.63	0.59	0.52	0.52	0.52	0.46	0.54	0.50	0.42	0.43	0.42	0.17	0.30	0.22
Micro-net	0.78	0.82	0.80	0.59	0.66	0.62	0.63	0.54	0.58	0.59	0.46	0.52	0.50	0.45	0.47	0.23	0.17	0.19
Hover-Net	0.82	0.79	0.80	0.58	0.67	0.62	0.54	0.60	0.56	0.56	0.51	0.54	0.52	0.47	0.49	0.28	0.35	0.31

- [9] K. Nagpal, D. Foote, Y. Liu, P.-H. C. Chen, E. Wulczyn, F. Tan, N. Olson, J. L. Smith, A. Mohtashamian, J. H. Wren, G. S. Corrado, R. MacDonald, L. H. Peng, M. B. Amin, A. J. Evans, A. R. Sangoi, C. H. Mermel, J. D. Hipp, and M. C. Stumpe, "Development and validation of a deep learning algorithm for improving Gleason scoring of prostate cancer," *npj Digital Medicine*, vol. 2, pp. 1–10, June 2019.
- [10] R. Colling, H. Pitman, K. Oien, N. Rajpoot, P. Macklin, D. Snead, T. Sackville, C. Verrill, C.-P. A. in Histopathology Working Group, V. Bachtiar, *et al.*, "Artificial intelligence in digital pathology: A roadmap to routine use in clinical practice," *The Journal of pathology*, 2019.
- [11] M. Shaban, S. A. Khurram, M. Hassan, S. Mushtaq, A. Loya, and N. Rajpoot, "Prognostic significance of automated score of tumor infiltrating lymphocytes in oral cancer," *Journal of Clinical Oncology*, vol. 36, pp. e18036–e18036, May 2018.
- [12] S. Javed, M. M. Fraz, D. Epstein, D. Snead, and N. M. Rajpoot, "Cellular Community Detection for Tissue Phenotyping in Histology Images," in *Computational Pathology and Ophthalmic Medical Image Analysis* (D. Stoyanov, Z. Taylor, F. Ciompi, Y. Xu, A. Martel, L. Maier-Hein, N. Rajpoot, J. van der Laak, M. Veta, S. McKenna, D. Snead, E. Trucco, M. K. Garvin, X. J. Chen, and H. Bogunovic, eds.), Lecture Notes in Computer Science, pp. 120–129, Springer International Publishing, 2018.
- [13] Y. Fu, A. W. Jung, R. V. Torne, S. Gonzalez, H. Vohringer, M. Jimenez-Linan, L. Moore, and M. Gerstung, "Pan-cancer computational histopathology reveals mutations, tumor composition and prognosis," *bioRxiv*, p. 813543, Oct. 2019.
- [14] Y. Saltz, R. Gupta, L. Hou, T. Kurc, P. Singh, V. Nguyen, D. Samaras, K. R. Shroyer, T. Zhao, B. E., Z. J. Heins, and Kundra, "Spatial Organization and Molecular Correlation of Tumor-Infiltrating Lymphocytes Using Deep Learning on Pathology Images," *Cell Reports*, vol. 23, pp. 181–193.e7, Apr. 2018.
- [15] K. Sirinukunwattana, E. Domingo, S. Richman, K. L. Redmond, A. Blake, C. Verrill, S. J. Leedham, A. Chatzipli, C. Hardy, C. Whalley, C.-H. Wu, A. D. Beggs, U. McDermott, P. Dunne, A. A. Meade, S. M. Walker, G. I. Murray, L. M. Samuel, M. Seymour, I. Tomlinson, P. Quirke, T. Maughan, J. Rittscher, V. H. Koelzer, and o. b. o. S. Consortium, "Image-based consensus molecular subtype classification (imCMS) of colorectal cancer using deep learning," *bioRxiv*, p. 645143, May 2019.
- [16] A. Reinke, S. Onogur, M. Stankovic, P. Scholz, T. Arbel, H. Bogunovic, A. P. Bradley, A. Carass, C. Feldmann, A. F. Frangi, *et al.*, "Is the winner really the best? a critical analysis of common research practice in biomedical image analysis competitions," *arXiv preprint arXiv:1806.02051*, 2018.
- [17] W. Brendel and M. Bethge, "Approximating CNNs with Bag-of-local-Features models works surprisingly well on ImageNet," Sept. 2018.
- [18] R. Geirhos, C. R. M. Temme, J. Rauber, H. H. Schtt, M. Bethge, and F. A. Wichmann, "Generalisation in humans and deep neural networks," in *Advances in Neural Information Processing Systems 31* (S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, eds.), pp. 7538–7550, Curran Associates, Inc., 2018.
- [19] J. Jo and Y. Bengio, "Measuring the tendency of CNNs to Learn Surface Statistical Regularities," *arXiv:1711.11561 [cs, stat]*, Nov. 2017. arXiv: 1711.11561.
- [20] J. Baxter, "A Model of Inductive Bias Learning," *Journal of Artificial Intelligence Research*, vol. 12, pp. 149–198, Mar. 2000. arXiv: 1106.0245.
- [21] J. Gamper, N. A. Koohbanani, K. Benet, A. Khuram, and N. Rajpoot, "PanNuke: An Open Pan-Cancer Histology Dataset for Nuclei Instance Segmentation and Classification," in *Digital Pathology*, pp. 11–19, Springer, Cham, Apr. 2019.
- [22] L. Oakden-Rayner, "Exploring large scale public medical image datasets," *arXiv:1907.12720 [cs, eess]*, July 2019. arXiv: 1907.12720.
- [23] S. Yune, H. Lee, S. Pomerantz, J. M. Romero, and S. Kamalian, "Real-world performance of deep-learning-based automated detection system for intracranial hemorrhages," in *SIIM Conference on Machine Intelligence and Medical Imaging*, 2018.
- [24] M. Jahanifar, N. A. Koohbanani, and N. Rajpoot, "NuClick: From Clicks in the Nuclei to Nuclear Boundaries," *arXiv:1909.03253 [cs, eess, q-bio, stat]*, Sept. 2019. arXiv: 1909.03253.
- [25] N. Kumar, R. Verma, S. Sharma, S. Bhargava, A. Vahadane, and A. Sethi, "A dataset and a technique for generalized nuclear segmentation for computational pathology," *IEEE transactions on medical imaging*, vol. 36, no. 7, pp. 1550–1560, 2017.
- [26] A. H. Beck, A. R. Sangoi, S. Leung, R. J. Marinelli, T. O. Nielsen, M. J. van de Vijver, R. B. West, M. van de Rijn, and D. Koller, "Systematic analysis of breast cancer morphology uncovers stromal features associated with survival," *Science Translational Medicine*, vol. 3, p. 108ra113, Nov. 2011.
- [27] H. Chang, J. Han, A. Borowsky, L. Loss, J. W. Gray, P. T. Spellman, and B. Parvin, "Invariant delineation of nuclear architecture in glioblastoma multiforme for clinical and molecular association," *IEEE transactions on medical imaging*, vol. 32, pp. 670–682, Apr. 2013.
- [28] A. Sethi, L. Sha, R. J. Deaton, V. Macias, A. H. Beck, and P. H. Gann, "Abstract LB-285: Computational pathology for predicting prostate cancer recurrence," *Cancer Research*, vol. 75, pp. LB–285–LB–285, Aug. 2015.
- [29] G. Lee, R. W. Veltri, G. Zhu, S. Ali, J. I. Epstein, and A. Madabhushi, "Nuclear Shape and Architecture in Benign Fields Predict Biochemical Recurrence in Prostate Cancer Patients Following Radical Prostatectomy: Preliminary Findings," *European Urology Focus*, vol. 3, no. 4–5, pp. 457–466, 2017.
- [30] M. S. Hosseini, L. Chan, G. Tse, M. Tang, J. Deng, S. Norouzi, C. Rowsell, K. N. Plataniotis, and S. Damaskinos, "Atlas of Digital Pathology: A Generalized Hierarchical Histological Tissue Type-Annotated Database for Deep Learning," p. 10.
- [31] L. Chan, M. S. Hosseini, C. Rowsell, K. N. Plataniotis, and S. Damaskinos, "Histosegnet: Semantic segmentation of histological tissue type in whole slide images," in *Proceedings of the IEEE International Conference on Computer Vision*, pp. 10662–10671, 2019.
- [32] K. Sirinukunwattana, J. P. Plum, H. Chen, X. Qi, P.-A. Heng, Y. B. Guo, L. Y. Wang, B. J. Matuszewski, E. Bruni, U. Sanchez, *et al.*, "Gland segmentation in colon histology images: The glas challenge contest," *Medical image analysis*, vol. 35, pp. 489–502, 2017.
- [33] Y. Fu, A. W. Jung, R. V. Torne, S. Gonzalez, H. Vohringer, M. Jimenez-Linan, L. Moore, and M. Gerstung, "Pan-cancer computational histopathology reveals mutations, tumor composition and prognosis," *bioRxiv*, p. 813543, 2019.
- [34] J. N. Kather, L. R. Heij, H. I. Grabsch, L. F. Kooreman, C. Loeffler, A. Echle, J. Krause, H. S. Muti, J. M. Niehues, K. A. Sommer, *et al.*, "Pan-cancer image-based detection of clinically actionable genetic alterations," *bioRxiv*, p. 833756, 2019.
- [35] K. Sirinukunwattana, S. E. A. Raza, Y. W. Tsang, D. R. J. Snead, I. A. Cree, and N. M. Rajpoot, "Locality Sensitive Deep Learning for Detection and Classification of Nuclei in Routine Colon Cancer Histology Images," *IEEE Transactions on Medical Imaging*, vol. 35, pp. 1196–1206, May 2016.
- [36] Q. D. Vu, S. Graham, T. Kurc, M. N. N. To, M. Shaban, T. Qaiser, N. A. Koohbanani, S. A. Khurram, J. Kalpathy-Cramer, T. Zhao, *et al.*, "Methods for segmentation and classification of digital microscopy tissue images," *Frontiers in bioengineering and biotechnology*, vol. 7, 2019.

- [37] J. Liu, T. Lichtenberg, K. A. Hoadley, L. M. Poisson, A. J. Lazar, A. D. Cherniack, A. J. Kovatich, C. C. Benz, D. A. Levine, A. V. Lee, *et al.*, “An integrated tcga pan-cancer clinical data resource to drive high-quality survival outcome analytics,” *Cell*, vol. 173, no. 2, pp. 400–416, 2018.
- [38] P. Kainz, M. Urschler, S. Schultze, P. Wohlhart, and V. Lepetit, “You Should Use Regression to Detect Cells,” in *Medical Image Computing and Computer-Assisted Intervention MICCAI 2015* (N. Navab, J. Hornegger, W. M. Wells, and A. F. Frangi, eds.), Lecture Notes in Computer Science, (Cham), pp. 276–283, Springer International Publishing, 2015.
- [39] L. Oakden-Rayner, J. Dunnmon, G. Carneiro, and C. R., “Hidden Stratification Causes Clinically Meaningful Failures in Machine Learning for Medical Imaging,” *arXiv:1909.12475 [cs, stat]*, Sept. 2019. arXiv: 1909.12475.
- [40] H. Qureshi, O. Sertel, N. Rajpoot, R. Wilson, and M. Gurcan, “Adaptive discriminant wavelet packet transform and local binary patterns for meningioma subtype classification,” *Medical image computing and computer-assisted intervention: MICCAI ... International Conference on Medical Image Computing and Computer-Assisted Intervention*, vol. 11, no. Pt 2, pp. 196–204, 2008.
- [41] P. Naylor, M. Laé, F. Reyal, and T. Walter, “Segmentation of nuclei in histopathology images by deep regression of the distance map,” *IEEE transactions on medical imaging*, vol. 38, no. 2, pp. 448–459, 2018.
- [42] K. He, G. Gkioxari, P. Dollár, and R. Girshick, “Mask r-cnn,” in *Proceedings of the IEEE international conference on computer vision*, pp. 2961–2969, 2017.
- [43] S. E. A. Raza, L. Cheung, M. Shaban, S. Graham, D. Epstein, S. Pelenaris, M. Khan, and N. M. Rajpoot, “Micro-net: A unified model for segmentation of various objects in microscopy images,” *Medical image analysis*, vol. 52, pp. 160–173, 2019.
- [44] Y. Xie, F. Xing, X. Shi, X. Kong, H. Su, and L. Yang, “Efficient and robust cell detection: A structured regression approach,” *Medical image analysis*, vol. 44, pp. 245–254, 2018.
- [45] A. Kirillov, K. He, R. Girshick, C. Rother, and P. Dollár, “Panoptic segmentation,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 9404–9413, 2019.
- [46] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, “ImageNet: A Large-Scale Hierarchical Image Database,” in *CVPR09*, 2009.
- [47] R. Werpachowski, A. Gyrgy, and C. Szepesvri, “Detecting Overfitting via Adversarial Examples,” *arXiv:1903.02380 [cs, stat]*, Mar. 2019. arXiv: 1903.02380.
- [48] A. Esteva, B. Kuprel, R. A. Novoa, J. Ko, S. M. Swetter, H. M. Blau, and S. Thrun, “Dermatologist-level classification of skin cancer with deep neural networks,” *Nature*, vol. 542, pp. 115–118, Feb. 2017.

APPENDIX

A. PanNuke labelling schema

The categories in PanNuke consist of: neoplastic, non-neoplastic epithelial, connective tissue, inflammatory and dead cells. These labels can be grouped into either neoplastic or

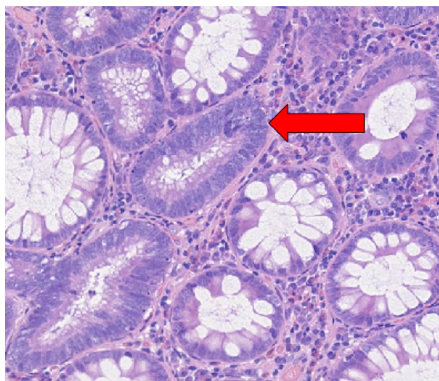


Fig. 9: Visual field extracted from colon tissue whole slide image. Red arrow points to a colorectal gland that consists of dysplastic epithelial cells.

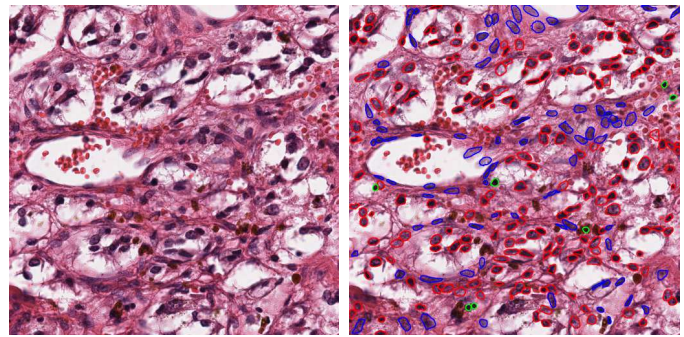


Fig. 10: Visual field extracted from adrenal gland tissue and its ground truth on the right.

non-neoplastic cell types. However, as Table A1 shows, non-neoplastic can cover everything from normal to inflammatory conditions, degenerative, meta-plastic, atypia and dysplasia. In Figure 9 we display a visual field from colon tissue where pathologists identified dysplastic epithelial cells. Here, dysplasia specifically refers to a pre-neoplastic stage, where it has not yet developed into a benign or malignant tumor. Therefore, these colon tissue cells would be labeled as non-neoplastic epithelial. Neoplastic labels in PanNuke specifically correspond to benign and malignant tumor cells.

TABLE A1: Breakdown of neoplastic and non-neoplastic cell types and sub-types.

Plasticity	Sub-types
Neoplastic	Malignant Tumors
	Benign Tumors
Non-neoplastic	Normal
	Hyperplastic
	Hypertrophic
	Meta-plastic
	Inflammatory
	Degenerative
	Atypia/Dysplasia

More than 2,000 visual fields from a variety of tissues have been reviewed when developing PanNuke. A major challenge for pathologists during the semi-automatic verification process in PanNuke was the necessity to refer back to the original WSIs to label the categories in the visual field. In Figure 10, parts of this image can be marked straightforwardly as stroma and inflammatory cells, but classifying the rest of this image is more challenging. For instance, it may represent retraction artifact or neoplasm, specifically pheochromocytoma. Once the WSI was viewed, this area turned out to be neoplastic, i.e. pheochromocytoma. This reiterates the point made in the main text that algorithms are trained on image patches and therefore do not contain the contextual information present in the WSI.

This is further exacerbated by the distribution of the nuclei size. Nuclei from the connective tissue are noticeably larger - however their shape and size do not directly indicate its category, as presented in Figure 11. This exemplifies the challenge of classifying between non-neoplastic epithelial and connective cell categories. Neoplastic cells tend to be far larger on average compared to any other category. Besides

the cell size variability per category within the tissue type, there is a significant difference in distribution of size between the tissues. Epithelial cells in particular vary in size between tissues- this emphasizes the importance of collecting labeled datasets across different tissue types.

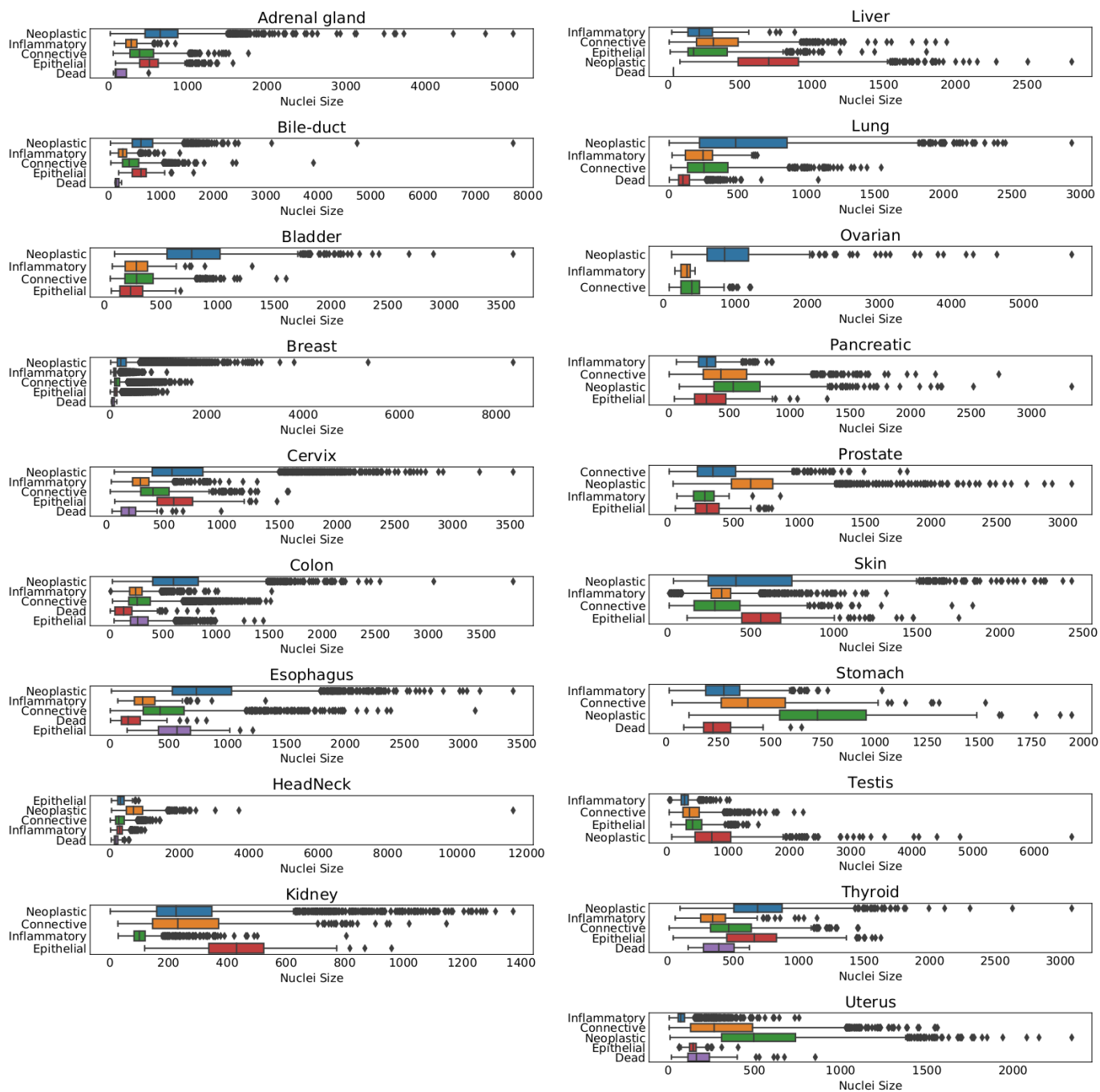


Fig. 11: Nuclei size distribution per class within every tissue type in PanNuke. Nuclei size is measured by pixel count within a segmentation mask for a particular nuclei.