

Variability Matters : Evaluating inter-rater variability in histopathology for robust cell detection

Cholmin Kang¹[0000-0001-6321-0003] Chunggi Lee¹[0000-0002-6164-2563], Heon Song¹[0000-0003-3435-0271], Minuk Ma¹[0000-0003-0416-8479], and Sérgio Pereira¹[0000-0002-4298-0903]

Lunit, Seoul 06241, Republic of Korea
 kcholmin@gmail.com, cglee@lunit.io, heon.song@lunit.io,
 akalsdnr3@gmail.com, sergio@lunit.io
<http://lunit.io>

Abstract. Large annotated datasets have been a key component in the success of deep learning. However, annotating medical images is challenging as it requires expertise and a large budget. In particular, annotating different types of cells in histopathology suffer from high inter- and intra-rater variability due to the ambiguity of the task. Under this setting, the relation between annotators’ variability and model performance has received little attention. We present a large-scale study on the variability of cell annotations among 120 board-certified pathologists and how it affects the performance of a deep learning model. We propose a method to measure such variability, and by excluding those annotators with low variability, we verify the trade-off between the amount of data and its quality. We found that naively increasing the data size at the expense of inter-rater variability does not necessarily lead to better-performing models in cell detection. Instead, decreasing the inter-rater variability with the expense of decreasing dataset size increased the model performance. Furthermore, models trained from data annotated with lower inter-labeler variability outperform those from higher inter-labeler variability. These findings suggest that the evaluation of the annotators may help tackle the fundamental budget issues in the histopathology domain

1 Introduction

Histopathology plays an important role in the diagnosis of cancer and its treatment planning. The process, typically, requires pathologists to localize cells and segment tissues in whole slide images (WSIs) [1]. However, with the advent of digital pathology, these tedious and error-prone tasks can be done by Deep Learning models. Still, these models are known to require the collection of large annotated datasets. Particularly, cell detection needs annotations from very large numbers of individually annotated cells.

The identification of cells in WSIs is challenging due to their diversity, subtle morphological differences, and the astounding amount of cells that exist in a

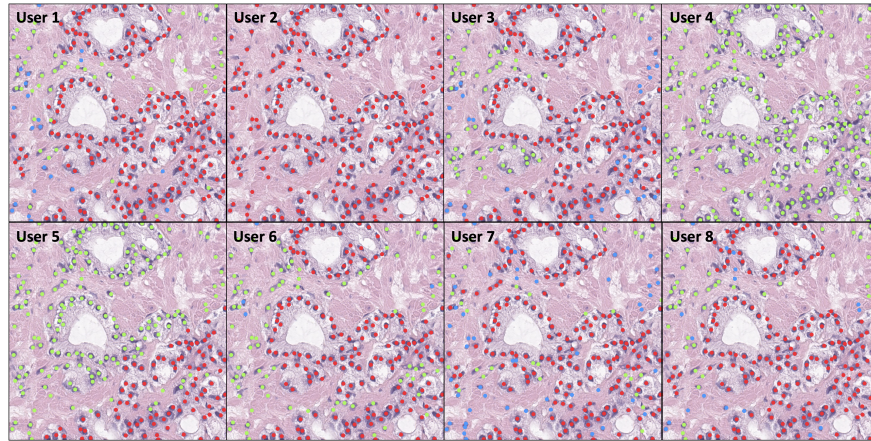


Fig. 1: Example of inter-rater variability in the annotation of cells. The colored circles represent: **red** - tumor cell, **blue** - lymphocytes, and **green** - other cells.

single WSI. Therefore, when manually performed, this task is time-consuming, and suffers from high inter- and intra-rater variability among annotators [2,3]. For example, Figure 1 shows the large discrepancy in cell annotations among different annotators.

The inconsistencies among the annotators negatively affect the performance of the model [4]. For example, when using such inconsistent data, the model may overfit to the noise and generalize poorly [5]. Therefore, some degree of consistency in the annotation must be guaranteed.

To improve the consistency among annotators and reduce the variability due to annotations in the dataset, we propose to exclude annotators with large variability from other annotators. This method is based on the existence of an anchor annotator who works as a reference point to which each annotator are compared. The concept of anchor has been recently proposed [6,7], where a cluster center is fixed for more efficient and stable training. We propose to have an anchor annotator that annotated a control set. In this way, we can measure the conformity of the other annotators in relation to the anchor one. During annotation, a given expert can receive a new image to annotate, or, instead, with some probability a control image. The set of annotated control images of an annotator is used to measure its conformity to the anchor annotator through a modified F1-score (mF1). In this way, annotators with high conformity would achieve lower variability, thus, more reliable annotations in relation to the anchor.

We collected and annotated a total of 29,387 patches from 12,326 WSIs. A total of 120 board-certified pathologists participated in the annotation process, which, to the best of our knowledge, is the largest-scale pathology annotation analysis reported in the literature. Based on our conformity, we divide the annotated data and perform an analysis on the discrepancy among annotators.

Our results show the conformity of 0.70, with a standard deviation of 0.08. Furthermore, the kappa score showed a mean agreement of $\kappa=0.43$ with a minimum value between $\kappa=-0.01$ and $\kappa=0.87$. These results show relatively low conformity among most of the annotators. Furthermore, we also verify that the data with lower variability leads to better-performing models. Indeed, the model trained with a smaller dataset with lower variability showed better performance when compared to the model trained with the full dataset. This supports the hypothesis that a smaller higher quality dataset can lead to a well-performing model, at a reduced cost, when compared with a large dataset with high inter-rater variability.

The contributions of this paper are five-fold. 1) We propose a method for calculating the conformity between annotators and an anchor annotator. In this way, we can measure the variability among annotators in cell detection tasks in histopathology. 2) We observe that the annotations created by pathologists have high inter-rater variability based on our conformity method. 3) By employing the before-mentioned conformity method, we propose a process to reduce annotation variability within a dataset by excluding annotators with low conformity in relation to an anchor. This leads to more efficient usage of the annotation budget, without sacrificing performance. 4) We show that a dataset with low inter-rater variability leads to better model performance. 5) To the best of our knowledge, this is the first work that tackles the problem of measuring variability among annotators in a large-scale dataset for a cell detection task, with the participation of a large crowd of board-certified pathologists.

2 Related work

Deep Learning techniques, particularly Convolutional Neural Networks (CNN) revolutionized the field of computer vision and image recognition [8,9] by achieving unprecedented performances. As such, it has also been extensively explored in the domain of digital pathology [10]. Cell detection tasks have received attention due to its laborious and error-prone nature [11,12,13,14]. However, there has been little attention to the effect of inter-rater variability in cell detection tasks, and how it affects the performance of the models.

In comparison to natural images, the annotation of medical images is a more challenging task. First, domain experts are needed, e.g., medical doctors. Second, due to the difficulty and ambiguity of the tasks, these tasks are more prone to subjectivity which leads to the disagreement between experts [15,16]. The inter-rater variability and reproducibility of annotations have been studied in some medical-related tasks, such as sleep pattern segmentation in EEG data [17], or segmentation in Computed Tomography images [18]. In the case of histopathology, some studies demonstrate that computer-assisted reading of images decreases the variability among observers [15,16].

To overcome the disagreements, some works attempt to utilize the annotations of multiple raters by learning several models from single annotations and from the agreement of multiple annotators [19]. In this way, the effect of dis-

agreement can be mitigated, but it assumes multiple annotations for the same images. Other approaches accept that inter-rater variability will exist, and explore crowds of non-experts to collect large amounts of data by crowd-sourcing, either with or without the assistance of experts. Nonetheless, the scopes were limited to tasks that do not require a high degree of expertise or the annotation of many instances [20,21]. Such assumptions do not hold in cell annotation in histopathology. Crowd-sourcing cell annotations were previously performed [22], but, the evaluation of the conformity of the annotators in this setting, where there is high inter-rater variability, remains unexplored. Inconsistencies among the annotators introduce noise to the annotations [4]. In turn, it can lead to models that overfit to the noise and generalize poorly [5]. Therefore, some degree of label variability must be guaranteed.

Other related prior work tackles the problem of lack of conformity between annotators from the perspective of a noisy label [23,24], or, tried to classify the reliability of annotators using the probabilistic model [25]. Other works tried to resolve large variability of annotator performance using the EM algorithm [26] or Gaussian process classifiers [27]. However, these approaches rely on the training dynamics, so, it is not possible to distinguish if a given sample is noisy or just difficult. Moreover, the noisy data is discarded, raising two concerns: 1) resources are used to gather unused data, and 2) the dataset needs to be sufficiently large to allow discarding data without impacting the performance of the model. Instead, we hypothesize that noise can be introduced by annotators with high variability to other annotators, and preventive selection of conforming annotators leads to higher quality annotations. This becomes even more necessary in large-scale settings, where a large pool of annotators needs to be sourced externally.

3 Evaluating annotators and variability

During the annotation process, the annotators receive unlabeled patches to annotate from the training set \mathbf{T} . Control set \mathbf{C} refers to the unlabeled dataset that will later be used to measure the variability among annotators. Images from the control set \mathbf{C} are randomly, and blindly, assigned to the annotators. Since the annotators do not know when a control image is provided, this ensures that the annotation skill showed in \mathbf{C} is similar to that of \mathbf{T} . Similar methods are used in practice when evaluating crowd-sourced workers [28]. An *anchor annotator* is chosen as the reference annotator to perform analysis on discrepancy (i.e., variability) among annotators. The annotations of each annotator are compared with the annotations done by the *anchor annotator*. Conformity is measured between each annotated cell of an annotator and the annotated cells of an *anchor annotator*. The conformity is used to show the variability of annotators and annotators with large conformity will have small variability among them.

Finally, annotators are divided into groups according to their conformity score. The model is, then, trained on the data of each group and aggregation of groups. The annotation, conformity measurement, and model development process are depicted in Figure 2(a).

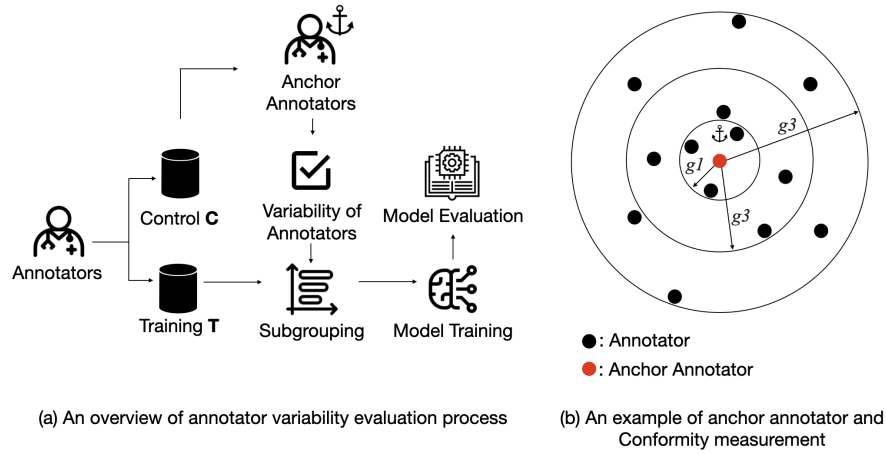


Fig. 2: Evaluation of inter-rater variability. (a) The control set is created to blindly evaluate the annotator conformity during the annotation process. A score is computed per annotator, which is used to evaluate the impact of their annotations in the cell detection task. (b) Annotator selection method. Conformity is measured between each annotation done by an annotator and anchor annotator. In the annotator selection phase, annotations are divided into subgroups (such as g_1 , g_2 , or g_3 in the figure) according to annotator conformity percentiles.

3.1 Anchor annotator

In order to perform an analysis of the variability among annotators, reference data should be set to measure how much each annotator is far from the reference point. In this work, one annotator is defined as the anchor annotator and is chosen as a criterion to evaluate variability among annotators. The rationale for setting an anchor annotator and evaluating variance is motivated by the concept of an anchor cluster center [6,7], where a cluster center is fixed for more efficient and stable training. The conceptual overview of anchor annotator selection is shown in Figure 3. As shown in the Figure, an anchor annotator is the one with most number of images annotated so that the annotated images are the superset of the patches annotated by other annotators.

3.2 Annotator variability and conformity calculation

Annotator conformity is calculated using algorithm 1. To compute the annotators' conformity, a hit criterion needs to be defined. We consider a circular region of radius $\theta = 10$ micrometers centered in each cell of \mathbf{C} in a given patch. Annotations from \mathbf{C} that are located inside a given circle are True Positive (TP) candidates. Among these candidates, the closest point that matches the same class as the annotation from \mathbf{C} is considered a TP. Given this criterion, we also compute False Positive (FP) and False Negative (FN) quantities. Based on

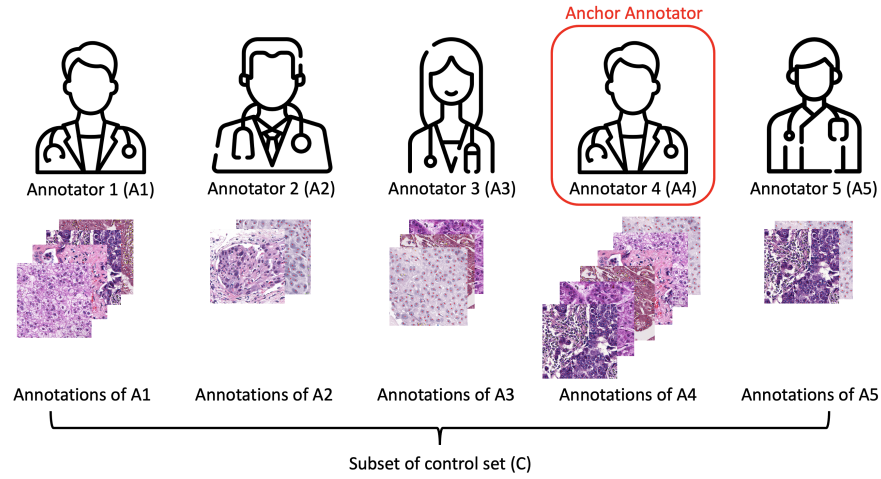


Fig. 3: A conceptual overview of anchor annotator selection. Anchor annotator is chosen among annotators with the most amount of annotations done in control set C . As shown in the figure, the patches annotated by A4 is the superset of patches annotated by A1, A2, A3, and A5. Thus, A4 is chosen as an anchor annotator.

the TP, FP, and FN, F1-score per class is computed and the conformity of the annotator is defined as mean F1-score over classes.

3.3 Annotator grouping based on the conformity

Finally, annotators are divided into subgroups according to their conformity to the anchor annotator. Figure 2(b) shows the annotators' annotation distribution. First, an anchor annotator is selected. Then, conformity between the anchor annotator and other annotators is calculated. Finally, each annotator is evaluated by computing the conformity, as described in algorithm 1. In this work, we consider this conformity as a variability. If all annotations done by all annotators are equal to those done by the anchor annotator, the conformity will be 1. However, if the annotations are different from the anchor annotator, the annotator conformity calculated by algorithm 1 will decrease, leading to increased variability and decreasing conformity.

4 Experimental Results

4.1 Dataset

We collected a large-scale dataset for the task of tumor and lymphocyte cells detection in H&E-stained WSIs, from 16 primary origins cancers (e.g. lung, breast, kidney, etc). As WSIs are typically large for exhaustive annotation, occupying

Algorithm 1: Conformity measurement algorithm for an individual annotator

Input: \mathbf{A} and \mathbf{C} : sets of annotations $\mathbf{A}_n^k = \{\mathbf{a}_{n,1}^k, \dots\}$ and $\mathbf{C}_n^k = \{\mathbf{c}_{n,1}^k, \dots\}$
 where $\mathbf{a}_{n,p}^k$ and $\mathbf{c}_{n,q}^k$ indicates the p -th and q -th annotations of the n -th control patch and class k

θ : distance threshold
Output: *conf*: conformity

$\mathbf{s} \leftarrow \mathbf{0}$; \triangleright Initialize a vector indicating F1-score for each class
foreach class k **do**
 $TP \leftarrow 0, FP \leftarrow 0$, and $FN \leftarrow 0$; \triangleright Initialize TP, FP, and FN
 foreach n -th control data patch **do**
 $G_A \leftarrow \emptyset$ and $G_M \leftarrow \emptyset$; \triangleright Initialize assigned index sets
 $D \in \mathbb{R}^{|\mathbf{A}_n^k| \times |\mathbf{C}_n^k|} \leftarrow \mathbf{0}$; \triangleright Initialize the pairwise distance matrix
 foreach index p and q **do**
 $D_{p,q} \leftarrow \|\mathbf{a}_{n,p}^k - \mathbf{c}_{n,q}^k\|_2$; \triangleright Compute the euclidean distance
 end
 foreach i -th index pair (p_i^*, q_i^*) such that $\forall j > i, D_{p_i^*, q_i^*} \leq D_{p_j^*, q_j^*}$ **do**
 if $D_{p_i^*, q_i^*} \leq \theta, p_i^* \notin G_A$, and $q_i^* \notin G_M$ **then**
 $G_A \leftarrow G_A \cup \{p_i^*\}$ and $G_M \leftarrow G_M \cup \{q_i^*\}$; \triangleright Assign indices
 end
 end
 $TP \leftarrow TP + |G_A|$; \triangleright Update TP
 $FP \leftarrow FP + |\mathbf{A}_n^k| - |G_A|$; \triangleright Update FP
 $FN \leftarrow FN + |\mathbf{C}_n^k| - |G_M|$; \triangleright Update FN
 end
 $s_k \leftarrow \frac{2TP}{2TP+FP+FN}$; \triangleright Compute the F1-score for class k
end
 $\text{conf} \leftarrow \frac{1}{|\mathbf{s}|} \sum_k s_k$; \triangleright Compute the conformity as mean F1-score

Table 1: Description of the Control, Training, and Validation sets in terms of the number patches, annotated cells, and WSIs.

Dataset	Patches	Annotated cells	Avg annotated cells per patch	Number of WSI
Control set	150	32,746	218.31	141
Training set	21,795	3,727,343	171.02	8,875
Validation set	7,442	1,390,551	186.85	3,310

several giga-pixels squared in area, it is a common practice to divide them into smaller patches [29]. We sampled patches of size 1024×1024 pixels from a total of 12,326 WSIs. The annotations consist of point annotations that locate the nuclei of cells, as well as their classes. We collect 3 sets at the WSI level: control set (**C**), Training set (**T**) and Validation set (**V**). The validation set is annotated by three annotators, and has no overlap with **C**. The number of images and cell

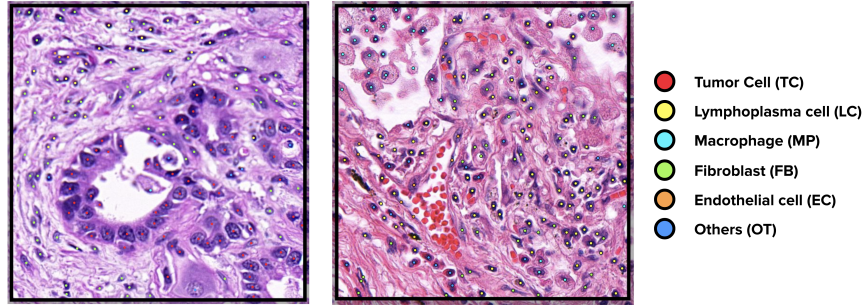


Fig. 4: Example of the reference images provided to annotators alongside the annotation guidelines.

annotations in each set is presented in Table 1. A total of 29,387 unique patches are annotated. As can be seen in the Table, the average number of annotated cells in a single image is similar among C , T , and V . A large crowd of 120 board-certified pathologists participated in the annotation process. During the annotation process, annotators are provided with guidelines on how to annotate each cell. Annotators are guided to annotate each cell to one of six classes: Tumor Cell, Lymphoplasma cells, Macrophage, Fibroblast, Endothelial cell, and others. The class "other" include nucleated cells with vague morphology that are not included in the specific cell types. Each cell is explained with its definition and its clinical characteristics. Furthermore, visual examples were provided to exemplify the guides, such as in Figure 4.

4.2 Cell detection method and experimental details

Having the conformity scores of the annotators, we evaluate its efficacy in filtering out annotations from annotators with higher variability by training a model in the downstream task of cell detection. Similar to [30], we pose the problem as a dense pixel classification task, but we use a DeepLabv3+ [31] architecture. The point annotation of each cell is converted to a circle of radius of 5 pixels, hence, resulting in a map similar to semantic segmentation tasks. The output is a dense prediction likelihood map that, similar to [30], requires a post-processing stage to retrieve the unique locations of the cells¹. In this experiment, we target Tumor-Infiltrating Lymphocyte (TIL) task which requires the number of tumor cells and lymphocytes in given WSI. Therefore, among the annotations, we regard macrophage, fibroblast, and endothelial cell as "others" class. All models are trained using the Adam [32] optimizer (learning rate of $1e-4$) and the soft dice loss function [33]. Experiments are performed on 4 NVIDIA V100 GPUs, and implemented using PyTorch (version 1.7.1) [34].

¹ Post-processing consists of Gaussian filtering ($\sigma = 3$) the likelihood maps, followed by local maxima detection within a radius of 3 pixels.

4.3 Conformity of the annotators

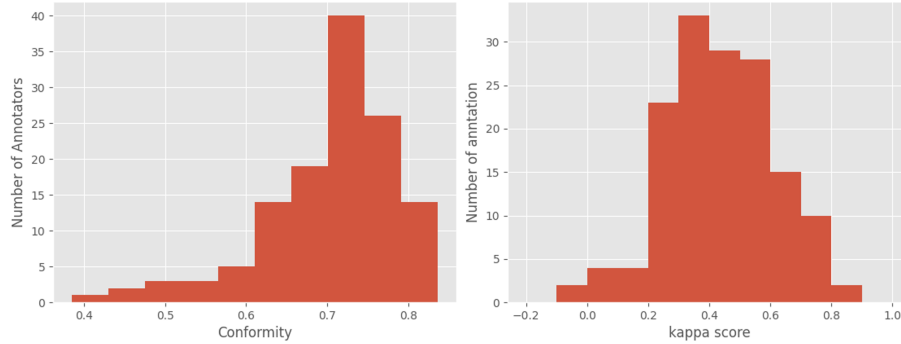


Fig. 5: Distributions of the variability of the external annotators as mF1-score (left), and the agreement of the annotations by the external annotators in relation to the set \mathbf{C} , measured in terms of Fleiss’ kappa (right).

Following the proposed evaluation procedure of the annotators, we computed the individual annotator conformity. In this work, an annotator with the most number of annotations done from set C is chosen as an anchor annotator so that most annotators can have conformity calculated. The closer a score is to 1, the better the annotator conforms with an anchor annotator. Figure 5, left, shows the distribution of the conformity score. The median conformity is 0.70 (minimum=0.38, and maximum=0.83). We further conduct a Shapiro-Wilk Normality test [35] and verify that it does not follow a Normal distribution at a significance level $\alpha = 0.05$, with a p-value= $9.427e - 8$. Moreover, we observe a negative skewness of -1.3158 , which confirms the observation of the longer left tail.

Varying expertise is a known problem when annotating cells in histopathology. Nonetheless, one would expect the conformity of the experts to reference annotation to be normally distributed. Instead, we observe that while the median annotator conformity is relatively high, there are some annotators on the left tail that perform significantly worse than the majority. This supports the observation of high inter-rater variability.

To further analyze the inter-rater variance, we measured the agreement among all annotators on each Control patch using Fleiss’ kappa [36], which is a statistical measure of the reliability of agreement for categorical ratings with more than two raters. This test was run to show inter-rater variability rather than to show absolute quality of annotators.

In this analysis, the inter-rater agreement is measured in relation to three categories: tumor cells, lymphocytes, and unmatched annotations. As previously described, the annotated cells from each pathologist have been matched to the nearest cell of the same class in \mathbf{C} , following the hit criterion. On the other hand,

Table 2: Model performance by each cancer primary origin and annotation conformity percentile range. For example, the column “ p_{0-25} ” denotes the training set with patches annotated by annotators with conformity in percentiles 0 to 25, i.e., the 25% annotators with the lowest conformity. At the bottom, it is presented the total annotation time in hours.

Primary Origin	Percentile ranges					
	p_{0-25}	p_{75-100}	p_{0-50}	p_{50-100}	p_{25-100}	p_{0-100}
Biliary Tract	62.09	62.9	63.71	64.06	65.1	64.89
Breast	62.32	64.74	64.53	65.53	64.83	65.3
Colorectum	61.87	63.09	63.83	63.8	64.02	63.84
Esophagus	64.11	61.42	65.56	66.24	66.33	66.19
Head and Neck	59.07	61.68	61.76	62.54	62.63	62.61
Kidney	49.68	53.94	52.10	56.07	57.34	56.49
Liver	62.50	65.16	64.06	66.09	65.29	65.16
Lung	58.18	60.24	60.27	61.24	60.75	60.72
Melanoma	59.08	60.03	60.18	61.64	61.73	61.27
Ovary	53.92	55.64	56.57	57.56	57.79	57.32
Pan Urinary	62.71	64.82	65.33	65.45	65.02	64.8
Pancreas	56.68	58.49	58.62	59.93	61.16	60.06
Prostate	55.70	56.41	56.12	58.13	58.39	58.64
Stomach	59.74	61.07	62.13	62.79	63.54	63.10
Uterine Cervix	61.86	64.33	64.37	65.23	65.69	65.64
Uterine Endometrium	53.18	54.70	54.13	56.07	55.53	55.51
Average mF1	58.92	60.77	60.87	62.03	62.18	61.76
Annotation time (hours)	598	575	1,122	1,176	1,703	2,298

if the hit criterion is not satisfied, we treat it as an unmatched annotation. We compute the Fleiss’ kappa for all patches in **C**. The mean agreement among the annotations of the pathologists was $\kappa=0.43$ with the minimum and maximum of $\kappa=-0.01$ and $\kappa=0.87$, respectively. This range (-0.01 to 0.87) of kappa values suggest poor to an excellent agreement, but the mean is located in a fair to moderate agreement range. This further reinforces the observation that there exists high inter-rater variability. We observe that the classification of the cells is sensitive in specific patches and pathologists, as shown in Figure 1.

4.4 Impact of the conformity of the annotators in cell detection

In this experiment, we divide the annotators by their conformity, from lower to higher. Following this conformity, we define 4 subsets of the data based on the percentile of the conformity of the annotators. Therefore, the interval p_{0-25} contains 25% of the patches in the dataset that were annotated by annotators with high variability. The intervals p_{25-50} , p_{50-75} , and p_{75-100} also contain 25% of the patches each, with decreasing annotator variability as compared to the anchor annotator. Therefore, following our hypothesis, the interval p_{75-100} must contain the annotations with the lowest variability. We hypothesize that data

annotated by annotators with lower variability should lead to better-performing cell detection models. Table 2 shows the performance of the cell detection model in terms of conformity by the before-mentioned intervals of data across the 16 cancer primary origin organs. The results are the average of 3 models trained with different random seeds. We use the paired Wilcoxon Signed Rank test [37] at a significance level $\alpha = 0.05$ to compare results and test our hypothesis.

We compare the performance of models trained with similar amounts of data, but with different conformity in relation to the annotation. First, we contrast interval p_{0-25} and p_{75-100} , i.e., the data annotated by the highest variance and lowest variance annotators, respectively. Note that both sets contain 25% of patches of the full training dataset. We observe that the data from the best performing annotators (interval p_{75-100}) leads to better performance compared to interval p_{0-25} . This is true for all of cancer primary origins and is corroborated by an existing statistical difference between the results (p-value=0.00614). Similarly, when we use 50% of the data to train the cell detection model (intervals p_{0-50} and p_{50-100}), a similar trend is observed. The model trained from the data from annotators with lower variability outperforms the model trained with data from interval p_{0-50} ; again, a statistical difference was found (p-value=0.00054). This result shows that given the same amount of data, the model trained with high conformity data outperforms that trained with low conformity data.

Additionally, we evaluate the performance of models trained from smaller datasets labeled by annotators with lower variability compared to the use of the full training set. The intervals p_{25-100} and p_{50-100} contain an upper 75% and 50% of the training dataset, respectively. Nonetheless, in both cases, they outperform the model trained with the full training set, with statistical significance in the case of the interval p_{25-100} (p-value=0.01778). These results suggest that smaller but low inter-rater variability datasets can outperform models trained from the larger datasets, but with higher variability annotations.

The data from the interval p_{50-100} represent the upper 50% of the available training data, but the performance is on par with the variability setting of the interval p_{25-100} . Therefore, in retrospect, we could have collected 50% fewer data, which would have reduced the total amount of annotation time from 2298 hours to 1176 hours. This is especially significant as the annotation process requires expert pathologists. Thus, we believe that the proposed annotation scoring procedure can be helpful in using the available budget more effectively, by requesting annotations from the experts with low variability to anchor annotator and obtain either a smaller dataset with good model performance or a large-scale dataset with lower inter-rater variability annotations in the future data collection process.

5 Conclusions

In this paper, we proposed a simple, yet effective method to assess the variability of a crowd of 120 experts in the time-consuming and difficult task of cell annotation in histopathology images stained with H&E. Our findings support

the previous observations that this task suffers from high inter-rater variability. Furthermore, we conducted a proof-of-concept experiment to show the effect of inter-rater variability on a machine learning model for the task of cell detection. From the experiment, we showed that with the same size of data, the data labeled by annotators with lower inter-rater variability led to better model performance than the data labeled by annotators with higher inter-rater variability. Furthermore, excluding the data with the most inter-rater variability from the full dataset was beneficial, despite resulting in a smaller dataset. Hence, contrary to common findings in deep learning research, the fewer amount of data with lower inter-rater variability resulted in better-performing machine learning model. These findings are useful in the medical domain where annotations are costly and laborious. With this finding, we can argue that collecting fewer data from low inter-rater variability is more beneficial than collecting data without considering the inter-rater variability. Despite the findings that increasing conformity leads to increasing model performance, we recognize that there are some limitations in our work. First, our work is based on the premise that the anchor annotator is well-performing. The work is based on the assumption that all annotators are sufficiently expert and have done their best to annotate. Further investigation is needed with iterated experiments on setting annotators as anchor annotators. Second, the generalizability of our work on stained WSI other than H&E stain is not guaranteed. Further investigation on other staining methods with more diverse skill backgrounds of digital pathologist annotators is needed. We hope the idea of scoring the annotators helps improve the annotation budget management in the histopathology domain.

References

1. Diao, J.A., Wang, J.K., Chui, W.F., Mountain, V., Gullapally, S.C., Srinivasan, R., Mitchell, R.N., Glass, B., Hoffman, S., Rao, S.K., et al.: Human-interpretable image features derived from densely mapped cancer pathology slides predict diverse molecular phenotypes. *Nature communications* **12**(1) (2021) 1–15
2. Lutnick, B., Ginley, B., Govind, D., McGarry, S.D., LaViolette, P.S., Yacoub, R., Jain, S., Tomaszewski, J.E., Jen, K.Y., Sarder, P.: An integrated iterative annotation technique for easing neural network training in medical image analysis. *Nature machine intelligence* **1**(2) (2019) 112–119
3. Cheplygina, V., Perez-Rovira, A., Kuo, W., Tiddens, H.A.W.M., de Bruijne, M.: Early experiences with crowdsourcing airway annotations in chest CT. In: MICCAI workshop. (2016)
4. Guo, B., Chen, H., Yu, Z., Xie, X., Huangfu, S., Zhang, D.: Fliermeet: a mobile crowdsensing system for cross-space public information reposting, tagging, and sharing. *Transactions on Mobile Computing* **14**(10) (2014) 2020–2033
5. Li, J., Wong, Y., Zhao, Q., Kankanhalli, M.S.: Learning to learn from noisy labeled data. In: CVPR. (2019)
6. Wang, Y., Wang, D., Pang, W., Miao, C., Tan, A.H., Zhou, Y.: A systematic density-based clustering method using anchor points. *Neurocomputing* **400** (2020) 352–370
7. Miller, D., Sunderhauf, N., Milford, M., Dayoub, F.: Class anchor clustering: A loss for distance-based open set recognition. In: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision. (2021) 3570–3578
8. Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems* **25** (2012)
9. Girshick, R., Donahue, J., Darrell, T., Malik, J.: Rich feature hierarchies for accurate object detection and semantic segmentation. In: Proceedings of the IEEE conference on computer vision and pattern recognition. (2014) 580–587
10. Meijering, E.: Cell segmentation: 50 years down the road [life sciences]. *IEEE Signal Processing Magazine* **29**(5) (2012) 140–145
11. Chen, T., Chef d’Hotel, C.: Deep learning based automatic immune cell detection for immunohistochemistry images. In: International workshop on machine learning in medical imaging, Springer (2014) 17–24
12. Gao, Z., Wang, L., Zhou, L., Zhang, J.: Hep-2 cell image classification with deep convolutional neural networks. *IEEE journal of biomedical and health informatics* **21**(2) (2016) 416–428
13. Cireşan, D.C., Giusti, A., Gambardella, L.M., Schmidhuber, J.: Mitosis detection in breast cancer histology images with deep neural networks. In: International conference on medical image computing and computer-assisted intervention, Springer (2013) 411–418
14. Xue, Y., Ray, N.: Cell detection in microscopy images with deep convolutional neural network and compressed sensing. *arXiv preprint arXiv:1708.03307* (2017)
15. Hekler, A., Utikal, J.S., Enk, A.H., Solass, W., Schmitt, M., Klode, J., Schandendorf, D., Sondermann, W., Franklin, C., Bestvater, F., et al.: Deep learning outperformed 11 pathologists in the classification of histopathological melanoma images. *European Journal of Cancer* **118** (2019) 91–96
16. Bertram, C.A., Aubreville, M., Donovan, T.A., Bartel, A., Wilm, F., Marzahl, C., Assenmacher, C.A., Becker, K., Bennett, M., Corner, S., et al.: Computer-assisted

- mitotic count using a deep learning-based algorithm improves interobserver reproducibility and accuracy. *Veterinary pathology* **59**(2) (2022) 211–226
17. Schaekermann, M., Beaton, G., Habib, M., Lim, A., Larson, K., Law, E.: Understanding expert disagreement in medical data analysis through structured adjudication. *Proceedings of the ACM on Human-Computer Interaction* **3**(CSCW) (2019) 1–23
 18. Haarbuerger, C., Müller-Franzes, G., Weninger, L., Kuhl, C., Truhn, D., Merhof, D.: Radiomics feature reproducibility under inter-rater variability in segmentations of ct images. *Scientific reports* **10**(1) (2020) 1–10
 19. Sudre, C.H., Anson, B.G., Ingala, S., Lane, C.D., Jimenez, D., Haider, L., Varsavsky, T., Tanno, R., Smith, L., Ourselin, S., et al.: Let’s agree to disagree: Learning highly debatable multirater labelling. In: *International Conference on Medical Image Computing and Computer-Assisted Intervention*, Springer (2019) 665–673
 20. Maier-Hein, L., Mersmann, S., Kondermann, D., Bodenstedt, S., Sanchez, A., Stock, C., Kenngott, H.G., Eisenmann, M., Speidel, S.: Can masses of non-experts train highly accurate image classifiers? - A crowdsourcing approach to instrument segmentation in laparoscopic images. In: *MICCAI*. (2014)
 21. Kwitt, R., Hegenbart, S., Rasiwasia, N., Vécsei, A., Uhl, A.: Do we need annotation experts? A case study in celiac disease classification. In: *MICCAI*. (2014)
 22. Amgad, M., Atteya, L.A., Hussein, H., Mohammed, K.H., Hafiz, E., Elsebaie, M.A., Alhusseiny, A.M., AlMoslemany, M.A., Elmatboly, A.M., Pappalardo, P.A., et al.: Nucls: A scalable crowdsourcing, deep learning approach and dataset for nucleus classification, localization and segmentation. *arXiv:2102.09099* (2021)
 23. Veit, A., Alldrin, N., Chechik, G., Krasin, I., Gupta, A., Belongie, S.: Learning from noisy large-scale datasets with minimal supervision. In: *CVPR*. (2017)
 24. Jiang, L., Zhou, Z., Leung, T., Li, L.J., Fei-Fei, L.: Mentornet: Learning data-driven curriculum for very deep neural networks on corrupted labels. In: *ICML*. (2018)
 25. Yan, Y., Rosales, R., Fung, G., Subramanian, R., Dy, J.: Learning from multiple annotators with varying expertise. *Machine learning* **95**(3) (2014) 291–327
 26. Raykar, V.C., Yu, S., Zhao, L.H., Valadez, G.H., Florin, C., Bogoni, L., Moy, L.: Learning from crowds. *Journal of machine learning research* **11**(4) (2010)
 27. Rodrigues, F., Pereira, F., Ribeiro, B.: Gaussian process classification and active learning with multiple annotators. In: *International conference on machine learning*, PMLR (2014) 433–441
 28. Bartolo, M., Thrush, T., Riedel, S., Stenetorp, P., Jia, R., Kiela, D.: Models in the loop: Aiding crowdworkers with generative annotation assistants. *arXiv:2112.09062* (2021)
 29. Khened, M., Kori, A., Rajkumar, H., Krishnamurthi, G., Srinivasan, B.: A generalized deep learning framework for whole-slide image segmentation and analysis. *Scientific reports* **11**(1) (2021) 1–14
 30. Swiderska-Chadaaj, Z., Pinckaers, H., van Rijthoven, M., Balkenhol, M., Melnikova, M., Geessink, O., Manson, Q., Sherman, M., Polonia, A., Parry, J., et al.: Learning to detect lymphocytes in immunohistochemistry with deep learning. *MIA* **58** (2019) 101547
 31. Chen, L., Zhu, Y., Papandreou, G., Schroff, F., Adam, H.: Encoder-decoder with atrous separable convolution for semantic image segmentation. In: *ECCV*. (2018)
 32. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. In Bengio, Y., LeCun, Y., eds.: *ICLR*. (2015)

33. Milletari, F., Navab, N., Ahmadi, S.A.: V-net: Fully convolutional neural networks for volumetric medical image segmentation. In: 3DV. (2016)
34. Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., Lin, Z., Desmaison, A., Antiga, L., Lerer, A.: Automatic differentiation in pytorch. (2017)
35. SHAPIRO, S.S., WILK, M.B.: An analysis of variance test for normality (complete samples)[†]. *Biometrika* **52**(3-4) (12 1965) 591–611
36. Falotico, R., Quatto, P.: Fleiss’ kappa statistic without paradoxes. *Quality & Quantity* **49**(2) (2015) 463–470
37. Woolson, R.F.: Wilcoxon signed-rank test. *Wiley encyclopedia of clinical trials* (2007) 1–3