

Title: Analytical validation of a standardized scoring protocol for Ki67 immunohistochemistry on breast cancer excision whole sections: an international multicenter collaboration.

Short running title: Standardized visual scoring of Ki67 in breast cancer

Authors: Samuel C.Y. Leung¹, Torsten O. Nielsen¹, Lila A. Zabaglo², Indu Arun³, Sunil S. Badve⁴, Anita L. Bane⁵, John M.S. Bartlett^{6,22}, Signe Borgquist⁷, Martin C. Chang⁸, Andrew Dodson⁹, Anna Ehinger¹⁰, Susan Fineberg¹¹, Cornelia M. Focke¹², Dongxia Gao¹, Allen M. Gown¹³, Carolina Gutierrez¹⁴, Judith C. Hugh¹⁵, Zuzana Kos¹⁶, Anne-Vibeke Lænkholm¹⁷, Mauro G. Mastropasqua¹⁸, Takuya Moriya¹⁹, Sharon Nofech-Mozes²⁰, C. Kent Osborne¹⁴, Frédérique M. Penault-Llorca²¹, Tammy Piper²², Takashi Sakatani²³, Roberto Salgado^{24,25}, Jane Starczynski²⁶, Tomoharu Sugie²⁷, Bert van der Vegt²⁸, Giuseppe Viale^{18,29}, Daniel F. Hayes³⁰, Lisa M. McShane³¹ and Mitch Dowsett² on behalf of the International Ki67 in Breast Cancer Working Group of the Breast International Group and North American Breast Cancer Group (BIG-NABCG)

Affiliations: ¹University of British Columbia, Vancouver, British Columbia, Canada; ²The Institute of Cancer Research, London, United Kingdom; ³Tata Medical Center, Kolkata, West Bengal, India; ⁴Indiana University Simon Cancer Center, Indianapolis, Indiana, United States; ⁵Juravinski Hospital and Cancer Centre, McMaster University, Hamilton, Ontario, Canada; ⁶Ontario Institute for Cancer Research, Toronto, Ontario, Canada; ⁷Division of Oncology and Pathology, Department of Clinical Science, Lund University, Lund, Sweden; ⁸Department of Pathology &

This is the author's manuscript of the article published in final edited form as:

Leung, S. C. Y., Nielsen, T. O., Zabaglo, L. A., Arun, I., Badve, S. S., Bane, A. L., Bartlett, J. M. S., Borgquist, S., Chang, M. C., Dodson, A., Ehinger, A., Fineberg, S., Focke, C. M., Gao, D., Gown, A. M., Gutierrez, C., Hugh, J. C., Kos, Z., Lænkholm, A.-V., ... Dowsett, M. (2019). Analytical validation of a standardised scoring protocol for Ki67 immunohistochemistry on breast cancer excision whole sections: An international multicentre collaboration. *Histopathology*, 75(2), 225–235. <https://doi.org/10.1111/his.13880>

Laboratory Medicine, University of Vermont Medical Center, Burlington, Vermont, United States; ⁹Ralph Lauren Centre for Breast Cancer Research, The Royal Marsden Hospital, London, United Kingdom; ¹⁰Department of Clinical Genetics and Pathology, Skane University Hospital, Lund University, Lund, Sweden; ¹¹Montefiore Medical Center and the Albert Einstein College of Medicine, Bronx, New York, United States; ¹²Dietrich-Bonhoeffer Medical Center, Neubrandenburg, Mecklenburg-Vorpommern, Germany; ¹³PhenoPath Laboratories, Seattle, Washington, United States; ¹⁴Lester and Sue Smith Breast Center and Dan L. Duncan Comprehensive Cancer Center, Baylor College of Medicine, Houston, Texas, United States; ¹⁵University of Alberta, Edmonton, Alberta, Canada; ¹⁶University of Ottawa and The Ottawa Hospital, Ottawa, Ontario, Canada; ¹⁷Department of Surgical Pathology, Zealand University Hospital, Slagelse, Region Sjælland, Denmark; ¹⁸European Institute of Oncology, Milan, Italy; ¹⁹Kawasaki Medical School, Kurashiki, Okayama Prefecture, Japan; ²⁰University of Toronto Sunnybrook Health Sciences Centre, Toronto, Ontario, Canada; ²¹Centre Jean Perrin and Université d'Auvergne, Clermont-Ferrand, France; ²²Edinburgh Cancer Research Centre, Western General Hospital, Edinburgh, United Kingdom; ²³Nippon Medical School, Bunkyo-ku, Tokyo, Japan; ²⁴Department of Pathology, GZA-ZNA, Antwerp, Belgium; ²⁵Division of Research, Peter MacCallum Cancer Centre, Melbourne, Australia; ²⁶Birmingham Heart of England, National Health Service, Birmingham, United Kingdom; ²⁷Kansai Medical University, Hirakata, Osaka, Japan; ²⁸University Medical Center Groningen, Groningen, Netherlands; ²⁹University of Milan, Milan, Italy; ³⁰University of Michigan Rogel Cancer Center, Ann Arbor, Michigan, United States and ³¹National Cancer Institute, Bethesda, Maryland, United States.

Corresponding author: Samuel Leung; Full postal address: Room 509, 2660 Oak Street, Jack Bell Research Center, Vancouver, BC, V6H 3Z6; Telephone: 604-875-4111 ext. 68893; Email: sam.leung@vch.ca

Acknowledgement of support: this work was supported by a generous grant from the Breast Cancer Research Foundation.

Conflict of interest statement:

Dr. Badve has participated in Scientific Advisory Boards/ Speaker for Genomic Health Inc., Dako/Agilent, Roche Diagnostics, Targos GmbH; Athenax, Konica-Minolta and received compensation. Dr. Badve has received research funding or in kind support from Dako/Agilent.

Dr. Badve has intellectual property right/ownership interests with IU. He is also associated with 2 startup companies (SYSGenomics and YeSSGenomics).

Dr. Bartlett has consulted for BioNTech GmbH, Biotheranostics Inc, RNA Diagnostics, and received compensation. Dr Bartlett has participated in Scientific Advisory Boards for Biotheranostics and RNA Diagnostics and received compensation. Dr Bartlett has received research funding or in kind support from Nanostring, Blotheranostics Inc, BioNTech GmbH. Dr Bartlett has intellectual property right/ownership interests with OICR/FACIT.

Dr. Borgquist has participated in educational talks / covered scientific conferences by Roche and Novartis.

Dr. Dowsett is on the Oncology Advisory Board for Radius and has provided ad hoc advice to Orion and Gtx. He has received lecture fees from Myriad and Roche and institutional research grants from Radius, Astrazeneca and Puma. He receives income from the Institute of Cancer Research Rewards for inventors Scheme (abiraterone).

Dr. Ehinger has participated in educational talks organized by Roche but without economical compensation.

Dr. Fineberg participated in a scientific advisory board for Genomic Health and have received monetary compensation (not for salary).

Dr. Hayes reports research support from Menarini Silicon Biosystems (MSB), Merrimack, Eli Lilly, Puma Biotechnology, Pfizer, AstraZeneca. He is the named inventor of patent US 8,790,878 B2. D.H.F. which is licensed to MSB and from whom he receives royalties. He holds stock options from Onclmmune LLC and InBiomotion, and he serves as a paid advisor for Cepheid, Freenome, CellWorks, Agendia, and CVS Caremark.

Dr. Lænkholm has received research funding from Nanostring Technology (not for personal salary), participated in advisory board for Roche A/S and Novartis (for purely scientific reasons; honoraria declined) and received travel expenses for congress attendance from Astra Zeneca and Roche A/S (past 2 years).

Dr. Nielsen has consulted for Nanostring and received compensation. Dr. Nielsen has intellectual property rights / ownership interests from Bioclassifier LLC.

Dr. Osborne has consulted for Astra Zeneca, Genentech and NanoString and received compensation.

Dr. Penault-Llorca has participated in Scientific Advisory Boards for Nanostring, Myriad, Genomic Health, Agendia, Astrazeneca, Roche, Sanofi, Novartis, Pfizer, BionTech, and received compensation. Dr. Penault-Llorca has received research funding or in kind support from Nanostring, Astrazeneca, Roche, BionTech.

Dr. Van der Vegt has consulted for Philips and received compensation.

All other authors declare no conflict of interest.

Word count: 2495

ABSTRACT

Aims: The nuclear proliferation marker Ki67 assayed by immunohistochemistry has multiple potential uses in breast cancer, but an unacceptable level of inter-laboratory variability has hampered its clinical utility. The International Ki67 in Breast Cancer Working Group has undertaken a systematic program to determine whether Ki67 measurement can be analytically validated and standardized across laboratories. This study addresses whether acceptable scoring reproducibility can be achieved on excision whole sections.

Methods and results: Adjacent sections from 30 primary ER+ breast cancers were centrally stained for Ki67 and sections were circulated among 23 pathologists in 12 countries. All pathologists scored Ki67 by two methods: (a) global: 4 fields of 100 tumor cells each were selected to reflect observed heterogeneity in nuclear staining; (b) hot-spot: the field with highest apparent Ki67 index was selected and up to 500 cells scored. The intraclass correlation coefficient (ICC) for the global method (0.87; 95%CI: 0.799-0.93) marginally met the prespecified success criterion (lower 95%CI \geq 0.8) while the ICC for the hot-spot method (0.83; 95%CI: 0.74-0.90) did not. Visually, inter-observer concordance in location of selected hot-spots varies between cases. The median times for scoring were 9 and 6 minutes for global and hot-spot methods, respectively.

Conclusions: The global scoring method demonstrates adequate reproducibility to warrant next steps toward evaluation for technical and clinical validity in appropriate cohorts of cases. The time taken for scoring by either method is practical using counting software we are making

publicly available. Establishment of external quality assessment schemes is likely to improve the reproducibility between laboratories further.

Keywords: Ki67, immunohistochemistry, pathology, scoring protocol, analytical validity, inter-observer variability, inter-observer reproducibility

INTRODUCTION

The nuclear antigen recognized by the Ki67 antibody is expressed in proliferating cells but absent in resting cells¹. Since its discovery in 1983 by Gerdes *et al.*,¹ Ki67 assessed by immunostaining has been studied extensively as a prognostic²⁻¹¹ and predictive^{4,6,9,12,13} marker, predominantly in hormone-receptor positive breast cancer but also in other tumors as well¹⁴⁻¹⁸. For example, pre-surgical Ki67 has been shown to be a marker for recurrence free survival¹⁹ and in the neoadjuvant setting, a marker for endocrine resistant tumor that may require more aggressive treatment²⁰. Excellent *intra*-observer reproducibility under controlled pre-analytic and staining conditions²¹ has contributed to the body of evidence showing the potential of Ki67 immunohistochemistry assay to be implemented in hospital laboratories as a cost effective part of clinical management²²⁻²⁴. However, poor inter-observer reproducibility and variability due to technical aspects of the assay has limited its adoption in clinical practice^{4,9,25-28}.

The International Ki67 Working Group (IKWG) has undertaken a systematic multiphase program to determine whether Ki67 *scoring* can be standardized and analytically validated across laboratories^{9,21,29,30}. In phase 1, as assessed by the intraclass correlation coefficient (ICC) estimate of inter-observer reproducibility, differences in pathologists' visual interpretation were the main source of variability (ICC = 0.71, 95% credible interval (CI): 0.47–0.78)²¹. In phase 2, greater concordance was achieved, at least on tissue microarrays, when pathologists trained to calibrate and standardize scoring according to a clearly defined methodology (ICC = 0.94, 95% CI: 0.90–0.97)²⁹. However, in clinical practice, decisions are made on core-cut biopsy or on excision specimens which require general assessment of the entire sample and selection of areas for formal counting. Therefore, in phase 3A, we assessed whether acceptable

performance could be achieved on core-cut biopsies using a standardized method with two distinct methods of scoring field selection: global (four representative fields, counting 100 nuclei each) and hot-spot (one field with highest Ki67, counting 500 nuclei). The global method achieved acceptable inter-observer reproducibility (ICC = 0.87; 95% CI: 0.81–0.93) according to our prespecified criteria, whereas the hot-spot method did not (ICC = 0.84; CI: 0.77–0.92)³⁰.

This current study represents the final phase (3B) of the visual scoring analytical validity program, wherein we assess whether acceptable performance can be achieved on centrally stained excision whole sections using the scoring method established on core-cut biopsies. Future studies would be required to evaluate variability due to staining and pre-analytical aspects of the assay.

MATERIALS AND METHODS

This study was approved by the British Columbia Cancer Agency Clinical Research Ethics Board (H10-03420). All specimens used in this study were donated by patients who signed institutionally-appropriate consent forms, were excess to diagnostic requirements and ethically available for quality control studies.

Case selection and sample preparation

Excision blocks from 30 estrogen receptor (ER) positive breast cancer cases were selected: 15 from the phase 3A study³⁰ and 15 from Kawasaki Medical School Hospital, Kurashiki, Japan (Supplemental Figure 1). Case selection was irrespective to patients' age at diagnosis, tumor grade, size or nodal status. The clinicopathological characteristics of these 30 cases are shown in Supplemental Table 1. All blocks were sectioned and stained in the Royal Marsden Hospital Histopathology Department using monoclonal antibody MIB1 at dilution 1:50 (DAKO UK, Cambridgeshire, UK) using an automated staining system (Ventana Medical Systems, Tucson, Arizona, USA) according to criteria established by the IKWG⁹. Sections from the same block were stained in a single immunohistochemistry run except for four cases where the staining was done in two different runs. This approach effectively controls for any technical variation in staining.

Sample distribution

Twenty-four volunteer pathologists representing 24 institutions from 12 countries, most of whom participated in the phase 3A study, were invited to participate.

Six adjacent sections from each of the 30 excision blocks were centrally stained: the first with H&E, the second with p63 (myoepithelial marker, to assist the identification of invasive foci) and the third to sixth with Ki67 (designated as slide sets 1–4). To facilitate application to the general histopathology laboratory environment, physical glass slides (as opposed to virtual slide images) were distributed to the volunteer pathologists. Because the accumulated delays required, if all pathologists reviewed the same physical glass slides, would have made the study impractical, participating pathologists were divided into four groups and were given one of the four sets of Ki67 slides to score. The H&E and p63 reference slides were made available online as digital images. Twenty-three pathologists successfully completed the study.

Scoring protocol

All pathologists were specifically trained to score Ki67 with emphasis on having a very low threshold for appreciating “brown stain” and the principles of standardized regions for nuclei counting, through the publicly available proficiency training module (<http://www.gpec.ubc.ca/calibrator>) that was initially used in the phase 2 study²⁹. The detailed scoring protocol is found in supplemental document: “ki67p3b_scoring_protocol.pdf”. A modified version of the scoring software used in this study is available freely from the Google Play and Apple iTunes store (search term: “Ki67”).

Scoring methods

The scoring methods used are the same ones that were employed in the phase 3A study³⁰: 1) a global assessment that is weighted according to the estimated percentage of the total cancer

area covered by each of high, medium, low, or negligible Ki67 staining levels; 2) an unweighted global assessment; and 3) assessment of Ki67 only in a “hot-spot” area.

Global methods attempt to derive an average score across all the tissue available for assessment. In the weighted and unweighted global methods, Ki67 index counting was performed in the same fashion, but the final Ki67 score was derived differently. Adapted from a scoring protocol that has been used routinely in the Dowsett laboratory³¹, these two global methods require the pathologist to first assess staining heterogeneity by estimating the percentages of the invasive tumor component of the slide exhibiting relatively high, medium, low or negligible Ki67 staining frequencies. Based on these estimates, an algorithm (Supplemental Figure 2) dictates the required number of fields to select and score for each Ki67 staining frequency (irrespective of staining intensity; totaling up to four fields). This algorithm was designed such that the four (or less) selected scoring fields would capture the full range of staining frequencies while at the same time, be reflective of the proportion in staining frequencies heterogeneity. Up to 100 invasive tumor nuclei within each field are counted using a “typewriter” pattern (Supplemental Figure 3), similar to how a tissue microarray core was scored in the phase 2 study²⁹.

The hot-spot method requires the pathologist to visually select one high-power field with the highest apparent staining rate and, within that area only, count up to 500 invasive tumor nuclei in a “typewriter” pattern.

Statistical analyses

Pre-specified criterion for success

Prior to data collection it was hypothesized that at least one of the scoring methods would have an associated ICC statistically greater than 0.80 (ICC of 0.8 being considered as good concordance³²). For planning purposes, power calculations performed under a variety of scenarios considered to represent good reproducibility (and similar to the results observed in the phase 2 study) showed that with at least 21 participating pathologists scoring 30 cases, there would be 80% power to exclude ICCs lower than the pre-specified ICC of 0.8 from a 95% credible interval for a given scoring method.

Ki67 score

The Ki67 score was defined as in the phase 3A study³⁰. Positive staining was defined as any brown stain in the nucleus above background, with reference available as needed to provide standard sample images; negative staining was scored when an invasive cancer cell showed only a blue counterstained nucleus. The unweighted global and hot-spot scores were simply the total number of positively stained tumor nuclei counted divided by the total number of tumor nuclei counted. The weighted global score was derived with tumor nuclei counts in each assessed field weighted by the estimated percentage of the total cancer area covered by each of high, medium, low, or negligible Ki67 staining levels. As in our previous studies, to satisfy model assumptions of normality and constant variance, for statistical analyses the Ki67 score is converted to a logarithmic scale by adding 0.1% and applying a log base 2 transformation.

ICC estimates (ranging from 0 to 1, with 1 representing perfect reproducibility) were computed as previously reported in the phase 3A study³⁰. Briefly, variance component analyses were performed to quantify the contributions from the following sources of variability: scoring

pathologist (observer), patient tumor (biological variation – each excision block represents a unique patient) and section of the excision block. Similar to the phase 3A study, same-section and different-section ICCs were computed. Same-section refers to pathologists scoring the same excision whole section physical slides, while different-section refers to pathologists scoring different physical slides that represent serial sections cut from the same original excision blocks. Credible intervals for the variance components and the ICCs were obtained using the Markov Chain Monte Carlo routines for fitting generalized linear mixed models.

All data analyses were performed using R version 3.3.2³³. Sources of variation in log₂-transformed Ki67 scores were analyzed using random effects models as implemented in the R packages lme4 and MCMCglmm. Data were visualized using heat maps, boxplots and spaghetti plots.

RESULTS

ICC of Ki67 according to scoring method.

The different-section ICC estimate for the weighted global scores was 0.87 (95%CI: 0.799–0.93), at the margin of the pre-specified success criterion (lower bound of credible interval exceeding 0.8) (Table 1). The different-section ICCs for the unweighted global scores and hot-spot scores were 0.86 (95%CI: 0.793–0.92) and 0.83 (95 %CI: 0.74–0.90), respectively, and therefore both these methods had ICC credible intervals that extended below the success criterion at the lower 95% limit. The corresponding same-section ICC estimates for the weighted global, unweighted global and hot-spot scores were virtually identical 0.87 (95% CI: 0.799–0.92), 0.86 (95% CI: 0.79–0.92) and 0.83 (95% CI: 0.74–0.90) respectively, supporting that differences between serial sections were minimal. Figure 1 displays the side-by-side boxplots of Ki67 scores across pathologists (hereafter referred to as “observers”) by group. Summary statistics for the Ki67 scores across the 23 observers are given in Supplemental Tables 2 to 4.

The median number of nuclei counted per slide (across all observers and cases) is 400 and 500 for the global and hot-spot methods respectively. The corresponding minimum number of nuclei counted is 300 and 138. Eighteen percent of the hot-spot scores were based on <500 nuclei counts. Among these 126 hot-spot scores, the median number of nuclei counted is 375.

In a context where preanalytical and staining factors are held constant, variance component analyses show that, regardless of scoring method, biological variation among different patients was the largest component of the total variation on these centrally stained slides, indicating

that the Ki67 score is reflecting inherent properties of the tumor (Figure 2, Supplemental Table 5).

Inter-observer variation of Ki67 scoring.

Figure 3 displays the variation in scores across observers for cases in slide set 1 as spaghetti plots. The corresponding plots for slide set 2-4 are displayed in Supplemental Figure 4. Figure 4 presents the scores in a heat map format with the columns (observers) ordered (within each slide set) by the median scores across cases and the rows (cases) sorted by the median scores across observers.

Overall it can be seen that most observers show good parallelism in the increasing Ki67 scores across the plots. In other words, observers measuring higher or lower than others tended to do so relatively consistently.

Categorical concordance of Ki67 scoring.

Regarding concordance on a categorical level (<10%, 10-20% and >20%), the relationship between concordance and continuous score is shown in Supplemental Figure 5. It shows excellent to perfect concordance on cases with scores that are either much lower or higher than the intermediate range (10-20%).

Based on visual inspection of captured images, locations of the hot-spot selections tended to cluster in the same region across observers within each of the excision whole section slides (Figure 5 shows some examples; virtual slide images of all slides used in this study and the

corresponding selected fields and scores can be viewed at

<http://www.gpec.ubc.ca/papers/ki67p3b>).

The median scoring time (field selection and nuclear counting) was 9 (interquartile range: 7-11) and 6 (interquartile range: 4-8) minutes for global and hot-spot methods, respectively.

DISCUSSION

The IKWG has demonstrated that it is possible, when controlling stringently for variability due to preanalytical and analytical aspects of the Ki67 immunohistochemistry assay⁹, and given a set of clearly defined training exercise and scoring instructions, for pathologists to achieve high inter-observer concordance in Ki67 scoring on core-cut biopsies and now on excision whole sections using a conventional light microscope and manual field selection, with no additional aid such as counting grid.

Due to the limited sample size, we were unable to assess whether any specific method (weighted global, unweighted global or hot-spot) is significantly more reproducible than others. However, the observed ICCs for global score (weighted: 0.87; unweighted: 0.86) are relatively higher compared to hot-spot score (0.83) suggesting that a sufficiently powered study might be able to show more convincingly whether global scores are more reproducible. This result is consistent with findings on core biopsies³⁰.

Can this level of concordance be clinically adequate? The POETIC¹¹ study assessed Ki67 (cut point at 10%) as a prognostic marker. Applying this cut point to the data in our current study, 17 (out of 30) cases have at most one discordance in weighted global score (Figure 4a). There are cases with major discrepancies: TB036, on the same physical slide (set 2), received a weighted global score of 4% and of 21% from observer A and L respectively. However, it is apparent (Figure 4) that cases far away from the intermediate range (10-20%) tend to have good agreement. Considering that cases in our current study are a random sampling of the general ER+ breast cancer population, one could expect that about half of these cases would

fall away from the intermediate range and hence Ki67 may provide clinically adequate information, provided that the staining and pre-analytical factors do not add too much variability.

Are the proposed scoring methods practical? The median scoring time is 6-9 minutes depending on the method used. However, an adaptive scoring protocol can be used to reduce scoring time if the purpose is to assess whether Ki67 is above or below a specific cut point. For example, considering the global scoring method, where the maximum nuclei count is pre-specified (i.e. 400), to determine whether a case has unweighted global score $\geq 10\%$, the pathologist can stop counting if the first field he/she scored is $\geq 40\%$. For cases with very low Ki67 score, one would likely still need to count all 400 nuclei.

The proposed scoring protocols do not make any recommendation concerning the required minimum tumor nuclei count. This is a limitation of this study and in practice, it will be up to the discretion of the scoring pathologist to assess if too few tumor nuclei are available for an adequate Ki67 assessment. This will depend on the percentage of positive cells scored in the cells available and the clinical context for the measurement.

External quality assessment program (e.g. NordiQC³⁴), involving comparing laboratory scores with reference scores in periodic assessment challenges, will likely improve inter-observer reproducibility further. Recent studies suggest that an even higher level of concordance can be achieved with automated image analysis³⁵⁻³⁸. The IKWG is actively conducting studies in this area to assess how artificial intelligence may help standardize Ki67 assessment^{35,38}. Also,

concordance between Ki67 scores on core biopsies and excision specimens is currently being investigated.

In conclusion, this study demonstrates an adequately high level of inter-observer concordance can be achieved by visual assessment of Ki67 using practical scoring methods, although some cases with large discrepancies remain. A two-tier assessment approach may be worthy of further study as a means to reduce scoring burden and further address challenging cases: if the Ki67 value from the initial scoring falls on a grey zone (e.g., cut point \pm 5%), scoring by a second pathologist or alternative test could be pursued. Preanalytical and analytical aspects of the immunohistochemistry assay, areas that still need standardization before the clinical utility of this marker can be proven, will likely add more variability. A clinical validation study employing analytically reproducible methodology would also need to be completed in appropriate cohorts of cases to determine whether Ki67 can be recommended for patient care decisions.

ACKNOWLEDGEMENTS

This work was supported by a generous grant from the Breast Cancer Research Foundation (DFH). Additional funding for the UK laboratories was received from Breakthrough Breast Cancer and the National Institute for Health Research Biomedical Research Centre at the Royal Marsden Hospital. Funding for work at the Ontario Institute for Cancer Research is provided by the Government of Ontario. Judith Hugh is the Lilian McCullough Chair in Breast Cancer Surgery Research and the CBCF Prairies/NWT Chapter. We are grateful to the Breast International Group and North American Breast Cancer Group (BIG-NABCG) collaboration, including the leadership of Drs. Nancy Davidson, Thomas Buchholz, Martine Piccart, and Larry Norton.

CONTRIBUTIONS

Samuel C.Y. Leung: study design, data collection, statistical analysis, manuscript drafting & review

Torsten O. Nielsen: study design, manuscript drafting & review

Lila A. Zabaglo: study design, collection and preparation of samples, data collection, manuscript drafting & review

Indu Arun: study design, data collection, manuscript drafting & review

Sunil S. Badve: study design, data collection, manuscript drafting & review

Anita L. Bane: study design, data collection, manuscript drafting & review

John M.S. Bartlett: study design, manuscript drafting & review

Signe Borgquist: study design, manuscript drafting & review

Martin C. Chang: study design, data collection, manuscript drafting & review

Andrew Dodson: study design, collection and preparation of samples, manuscript drafting & review

Anna Ehinger: study design, data collection, manuscript drafting & review

Susan Fineberg: study design, data collection, manuscript drafting & review

Cornelia M. Focke: study design, data collection, manuscript drafting & review

Dongxia Gao: study design, data collection, manuscript drafting & review

Allen M. Gown: study design, data collection, manuscript drafting & review

Carolina Gutierrez: study design, data collection, manuscript drafting & review

Judith C. Hugh: study design, data collection, manuscript drafting & review

Zuzana Kos: study design, data collection, manuscript drafting & review

Anne-Vibeke Lænkholm: study design, data collection, manuscript drafting & review

Mauro G. Mastropasqua: study design, data collection, manuscript drafting & review

Takuya Moriya: study design, data collection, manuscript drafting & review

Sharon Nofech-Mozes: study design, data collection, manuscript drafting & review

C. Kent Osborne: study design, manuscript drafting & review

Frédérique M. Penault-Llorca: study design, data collection, manuscript drafting & review

Tammy Piper: study design, data collection, manuscript drafting & review

Takashi Sakatani: study design, data collection, manuscript drafting & review

Roberto Salgado: study design, data collection, manuscript drafting & review

Jane Starczynski: study design, data collection, manuscript drafting & review

Tomoharu Sugie: study design, manuscript drafting & review

Bert van der Vegt: study design, data collection, manuscript drafting & review

Giuseppe Viale: study design, manuscript drafting & review

Daniel F. Hayes: study design, manuscript drafting & review

Lisa M. McShane: study design, statistical analysis, manuscript drafting & review

Mitch Dowsett: study design, manuscript drafting & review

FUNDING

This work was supported by a generous grant from the Breast Cancer Research Foundation.

REFERENCES

- (1) Gerdes J, Schwab U, Lemke H, Stein H. Production of a mouse monoclonal antibody reactive with a human nuclear antigen associated with cell proliferation. *Int J Cancer* 1983 Jan 15;31(1):13-20.
- (2) Luporsi E, Andre F, Spyrtos F et al. Ki-67: level of evidence and methodological considerations for its role in the clinical management of breast cancer: analytical and critical review. *Breast Cancer Res Treat* 2012 Apr;132(3):895-915.
- (3) de Azambuja E, Cardoso F, de Castro G, Jr et al. Ki-67 as prognostic marker in early breast cancer: a meta-analysis of published studies involving 12,155 patients. *Br J Cancer* 2007 May 21;96(10):1504-1513.
- (4) Denkert C, Budczies J, von Minckwitz G, Wienert S, Loibl S, Klauschen F. Strategies for developing Ki67 as a useful biomarker in breast cancer. *Breast* 2015 Aug 14.
- (5) Inwald EC, Klinkhammer-Schalke M, Hofstadter F et al. Ki-67 is a prognostic parameter in breast cancer patients: results of a large population-based cohort of a cancer registry. *Breast Cancer Res Treat* 2013 Jun;139(2):539-552.
- (6) Viale G, Regan MM, Maiorano E et al. Prognostic and predictive value of centrally reviewed expression of estrogen and progesterone receptors in a randomized trial comparing letrozole and tamoxifen adjuvant therapy for postmenopausal early breast cancer: BIG 1-98. *J Clin Oncol* 2007 Sep 1;25(25):3846-3852.

(7) Viale G, Regan MM, Mastropasqua MG et al. Predictive value of tumor Ki-67 expression in two randomized trials of adjuvant chemoendocrine therapy for node-negative breast cancer. *J Natl Cancer Inst* 2008 Feb 6;100(3):207-212.

(8) Yerushalmi R, Woods R, Ravdin PM, Hayes MM, Gelmon KA. Ki67 in breast cancer: prognostic and predictive potential. *Lancet Oncol* 2010 Feb;11(2):174-183.

(9) Dowsett M, Nielsen TO, A'Hern R et al. Assessment of Ki67 in breast cancer: recommendations from the International Ki67 in Breast Cancer working group. *J Natl Cancer Inst* 2011 Nov 16;103(22):1656-1664.

(10) Petrelli F, Viale G, Cabiddu M, Barni S. Prognostic value of different cut-off levels of Ki-67 in breast cancer: a systematic review and meta-analysis of 64,196 patients. *Breast Cancer Res Treat* 2015 Sep 4.

(11) Robertson JFR, Dowsett M, Bliss JM et al. Peri-operative aromatase inhibitor treatment in determining or predicting longterm outcome in early breast cancer - The POETIC Trial. San Antonio Breast Cancer Symposium 2017 presented December 6, 2017;abstract GS1-03.

(12) Criscitiello C, Disalvatore D, De Laurentiis M et al. High Ki-67 score is indicative of a greater benefit from adjuvant chemotherapy when added to endocrine therapy in luminal B HER2 negative and node-positive breast cancer. *Breast* 2014 Feb;23(1):69-75.

- (13) Cohen AL, Factor RE, Mooney K et al. POWERPIINC (PreOperative Window of Endocrine Therapy Provides Information to Increase Compliance) trial: Changes in tumor proliferation index and quality of life with 7 days of preoperative tamoxifen. *Breast* 2017 Feb;31:219-223.
- (14) Lei Y, Li Z, Qi L et al. The Prognostic Role of Ki-67/MIB-1 in Upper Urinary-Tract Urothelial Carcinomas: A Systematic Review and Meta-analysis. *J Endourol* 2015 Jul 23.
- (15) Desouki MM, Chamberlain BK, Li Z. The role of immunohistochemistry in the evaluation of gynecologic pathology part 2: a comparative study between two academic institutes. *Ann Diagn Pathol* 2015 Jun 11.
- (16) He Y, Wang N, Zhou X et al. Prognostic value of ki67 in BCG-treated non-muscle invasive bladder cancer: a meta-analysis and systematic review. *BMJ Open* 2018 Apr 17;8(4):e019635-2017-019635.
- (17) Richardsen E, Andersen S, Al-Saad S et al. Evaluation of the proliferation marker Ki-67 in a large prostatectomy cohort. *PLoS One* 2017 Nov 15;12(11):e0186852.
- (18) Xie Y, Chen L, Ma X et al. Prognostic and clinicopathological role of high Ki-67 expression in patients with renal cell carcinoma: a systematic review and meta-analysis. *Sci Rep* 2017 Mar 13;7:44281.
- (19) Dowsett M, Smith IE, Ebbs SR et al. Prognostic value of Ki67 expression after short-term presurgical endocrine therapy for primary breast cancer. *J Natl Cancer Inst* 2007 Jan 17;99(2):167-170.

- (20) Ellis MJ, Suman VJ, Hoog J et al. Ki67 Proliferation Index as a Tool for Chemotherapy Decisions During and After Neoadjuvant Aromatase Inhibitor Treatment of Breast Cancer: Results From the American College of Surgeons Oncology Group Z1031 Trial (Alliance). *J Clin Oncol* 2017 Apr 1;35(10):1061-1069.
- (21) Polley MY, Leung SC, McShane LM et al. An international Ki67 reproducibility study. *J Natl Cancer Inst* 2013 Dec 18;105(24):1897-1906.
- (22) Iwamoto T, Katagiri T, Niikura N et al. Immunohistochemical Ki67 after short-term hormone therapy identifies low-risk breast cancers as reliably as genomic markers. *Oncotarget* 2017 Apr 18;8(16):26122-26128.
- (23) Thakur SS, Li H, Chan AMY et al. The use of automated Ki67 analysis to predict Oncotype DX risk-of-recurrence categories in early-stage breast cancer. *PLoS One* 2018 Jan 5;13(1):e0188983.
- (24) Reinert T, Goncalves R, Ellis MJ. Current Status of Neoadjuvant Endocrine Therapy in Early Stage Breast Cancer. *Curr Treat Options Oncol* 2018 Apr 16;19(5):23-018-0538-9.
- (25) Laenkholm AV, Grabau D, Moller Talman ML et al. An inter-observer Ki67 reproducibility study applying two different assessment methods: on behalf of the Danish Scientific Committee of Pathology, Danish breast cancer cooperative group (DBCG). *Acta Oncol* 2018 Jan;57(1):83-89.
- (26) Focke CM, Burger H, van Diest PJ et al. Interlaboratory variability of Ki67 staining in breast cancer. *Eur J Cancer* 2017 Oct;84:219-227.

- (27) Mengel M, von Wasielewski R, Wiese B, Rudiger T, Muller-Hermelink HK, Kreipe H. Inter-laboratory and inter-observer reproducibility of immunohistochemical assessment of the Ki-67 labelling index in a large multi-centre trial. *J Pathol* 2002 Nov;198(3):292-299.
- (28) Ekholm M, Grabau D, Bendahl PO et al. Highly reproducible results of breast cancer biomarkers when analysed in accordance with national guidelines - a Swedish survey with central re-assessment. *Acta Oncol* 2015 Jul;54(7):1040-1048.
- (29) Polley MY, Leung SC, Gao D et al. An international study to increase concordance in Ki67 scoring. *Mod Pathol* 2015 Jun;28(6):778-786.
- (30) Leung SCY, Nielsen TO, Zabaglo L et al. Analytical validation of a standardized scoring protocol for Ki67: phase 3 of an international multicenter collaboration. *NPJ Breast Cancer* 2016 May 18;2:16014.
- (31) Zabaglo L, Salter J, Anderson H et al. Comparative validation of the SP6 antibody to Ki67 in breast cancer. *J Clin Pathol* 2010 Sep;63(9):800-804.
- (32) Kirkegaard T, Edwards J, Tovey S et al. Observer variation in immunohistochemical analysis of protein expression, time for a change? *Histopathology* 2006 Jun;48(7):787-794.
- (33) R Core Team. R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. 2018.
- (34) Vyberg M, Møller J, Røge R. Nordic immunohistochemical Quality Control - Ki67 assessment. 2018; Available at: <http://www.nordiqc.org/epitope.php?id=1>, 2018.

(35) Acs B, Pelekanou V, Bai Y et al. Ki67 reproducibility using digital image analysis: an inter-platform and inter-operator study. *Lab Invest* 2018 Sep 4.

(36) Stalhammar G, Robertson S, Wedlund L et al. Digital image analysis of Ki67 in hot spots is superior to both manual Ki67 and mitotic counts in breast cancer. *Histopathology* 2018 May;72(6):974-989.

(37) Koopman T, Buikema HJ, Hollema H, de Bock GH, van der Vegt B. Digital image analysis of Ki67 proliferation index in breast cancer using virtual dual staining on whole tissue sections: clinical validation and inter-platform agreement. *Breast Cancer Res Treat* 2018 May;169(1):33-42.

(38) Rimm DL, Leung SCY, McShane LM et al. An international multicenter study to evaluate reproducibility of automated scoring for assessment of Ki67 in breast cancer. *Mod Pathol* 2018 Aug 24.

TABLE**Table 1.** Summary of ICC values for different scoring methods.

	Different-section ICC	Same-section ICC
Weighted global	0.87 (95%CI: 0.799–0.93)	0.87 (95% CI: 0.799–0.92)
Unweighted global	0.86 (95%CI: 0.79–0.92)	0.86 (95% CI: 0.79–0.92)
Hot-spot	0.83 (95 %CI: 0.74–0.90)	0.83 (95% CI: 0.74–0.90)

FIGURE LEGENDS

Figure 1. Ki67 scores of all 23 observers (by slide set). Observers are ordered (within each group) by the median scores. The bottom/top of the box in each box plot represent the first (Q1)/third (Q3) quartiles, the bold line inside the box represents the median and the two bars outside the box represent the lowest/highest datum still within $1.5 \times$ the inter-quartile range (Q3-Q1). Outliers are represented with empty circles.

Figure 2. Variance component analysis. Variation due to different components are presented in a bar plot to show the relative magnitude of differences between them. Numeric values of the variance components estimates and the corresponding credible intervals are shown in Supplemental Table 5.

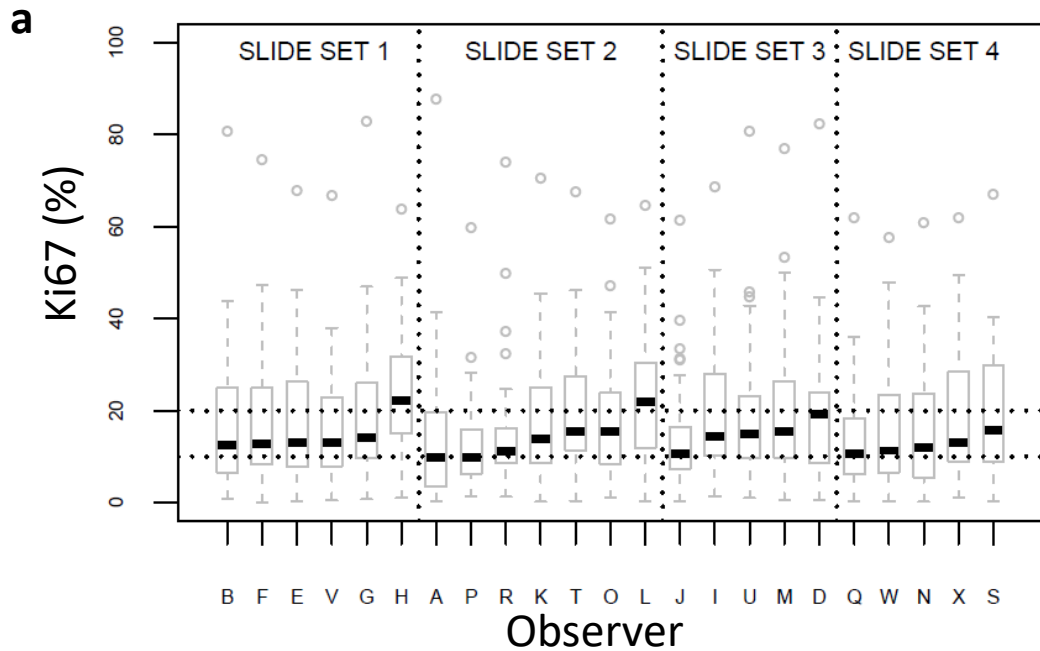
Figure 3. Variability in Ki67 scores (slide set 1 only). Each line represents Ki67 scores from one observer. Shaded region indicates Ki67 scores between 10-20%. Scores on slide set 2-4 are shown in Supplemental Figure 4.

Figure 4. Heat map of Ki67 scores (a: weighted global; b: unweighted global; c: hot-spot). Rows represent cases and columns represent observers. Green color indicates that the score is <10%, yellow 10-20% and red >20%. Cases are ordered by the median scores (across observers), which are shown in parentheses beside the specimen number. Observers are ordered (within each group) by the median scores (across cases). The three colon-separated numbers to the right of the heat map represent the number of observers giving scores falling into different ranges: <10% (left-most), 10-20% (middle) and >20% (right-most). For example, "15:6:1"

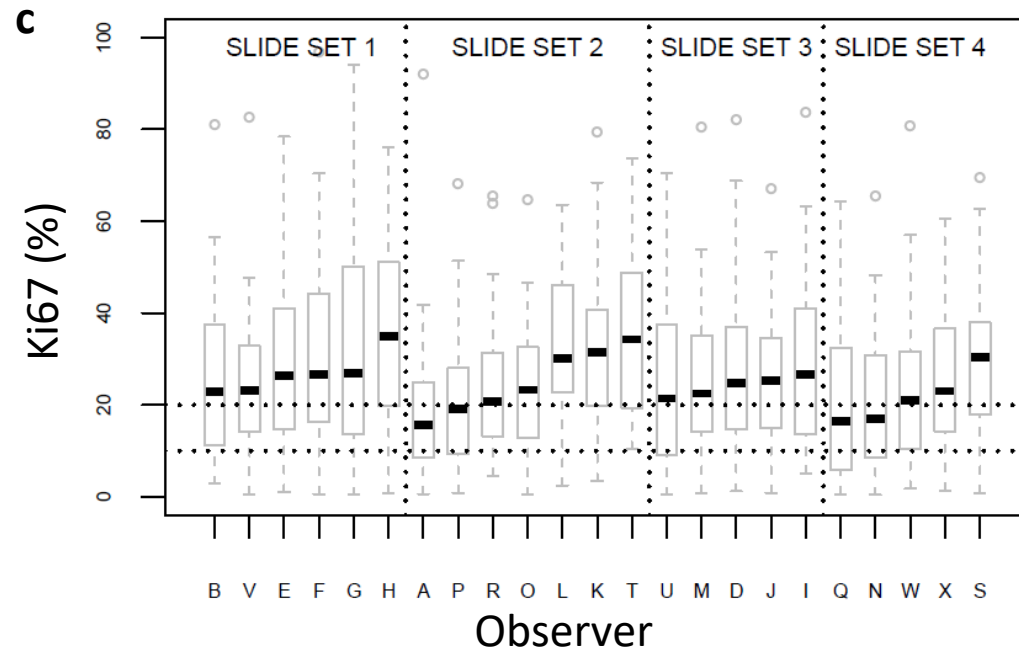
indicates that 15 observers gave a score of <10%, six observers between 10-20% and one observer >20%.

Figure 5. Hot-spot field selection by different observers on the same excision whole section slide. Figure 5a shows selections (indicated by red circles) on some example excision whole section slides. Figure 5b is an example of a single excision whole section slide (median score: 18%) with zoomed-in fields. Each observer was asked to circle the area considered to be the hot spot (b-i). Most observers honed in on the same general area of the slide, although individual selected scoring fields do not always overlap. Figure 5b-iii and 5b-iv represent segments of the same area chosen by two different observers to read Ki67. Figure 5b-v represents the “outlier” field selected by only one observer as the hot-spot.

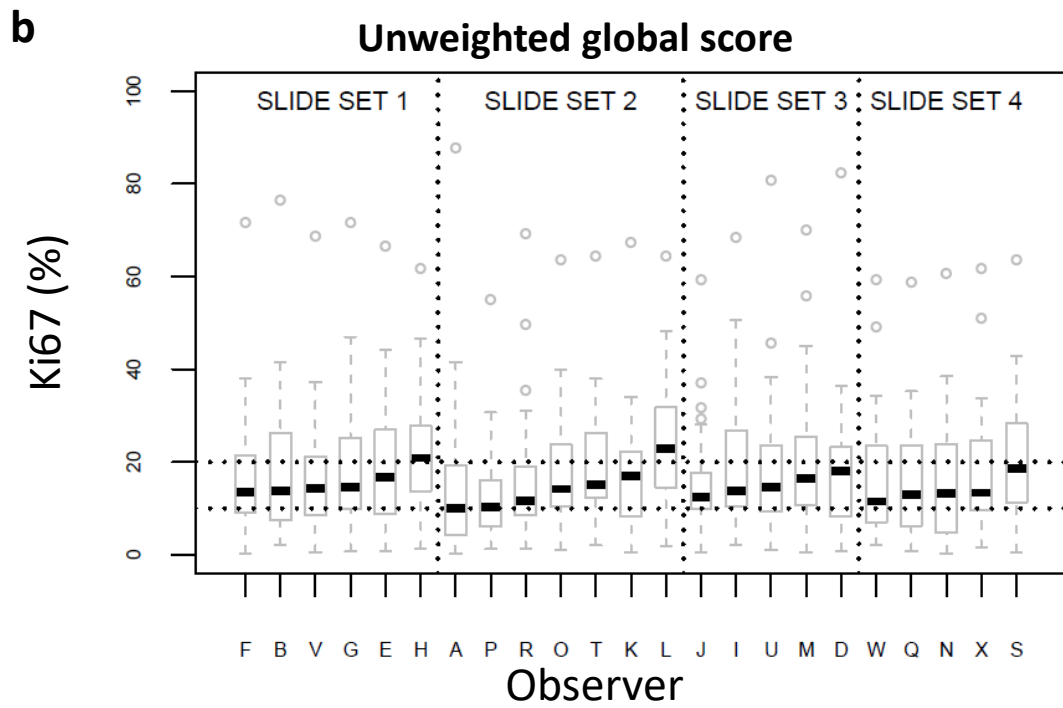
Weighted global score



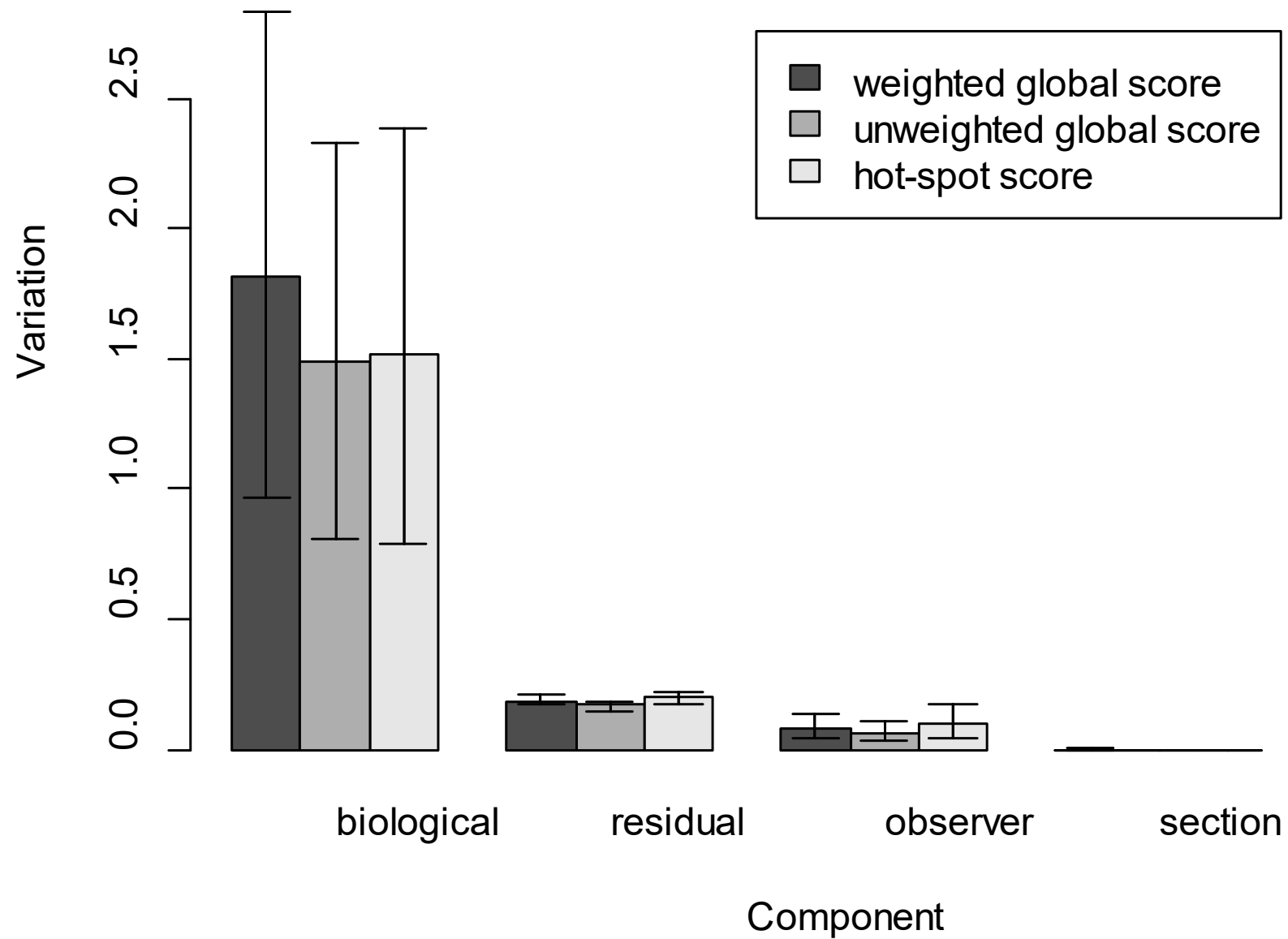
Hot-spot score

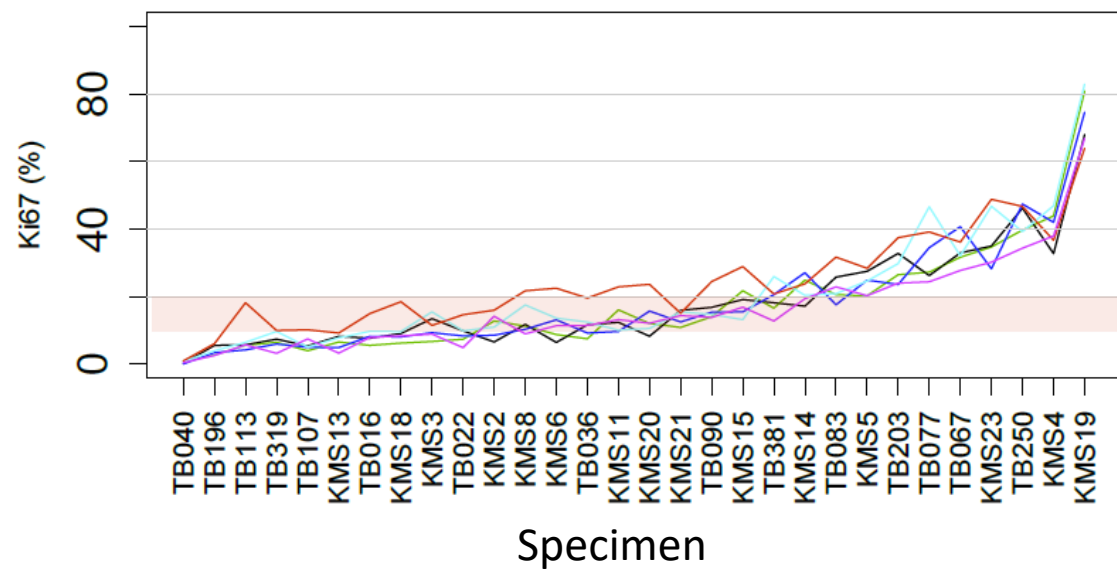
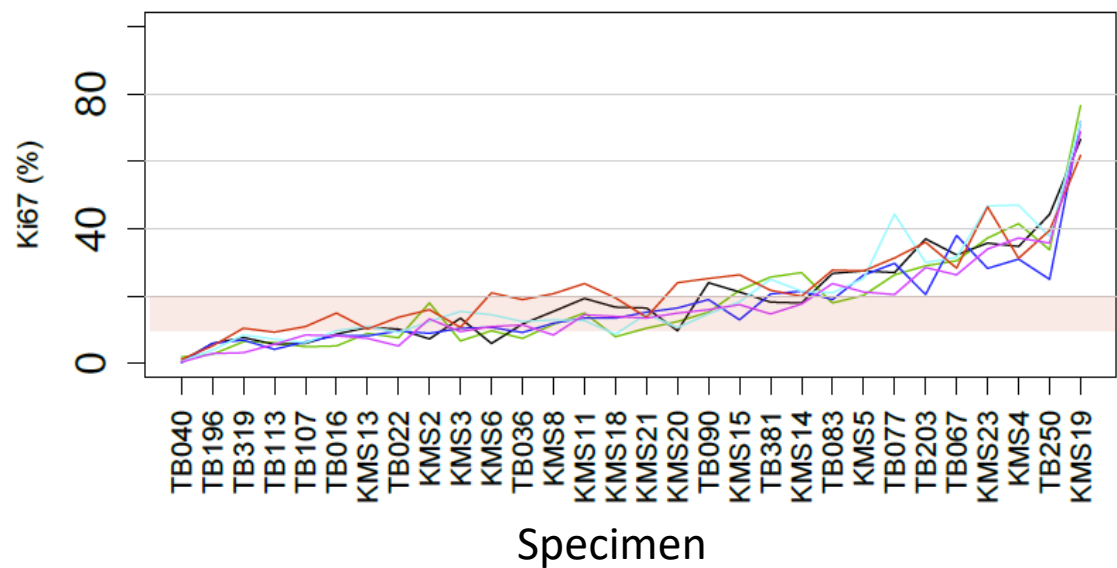
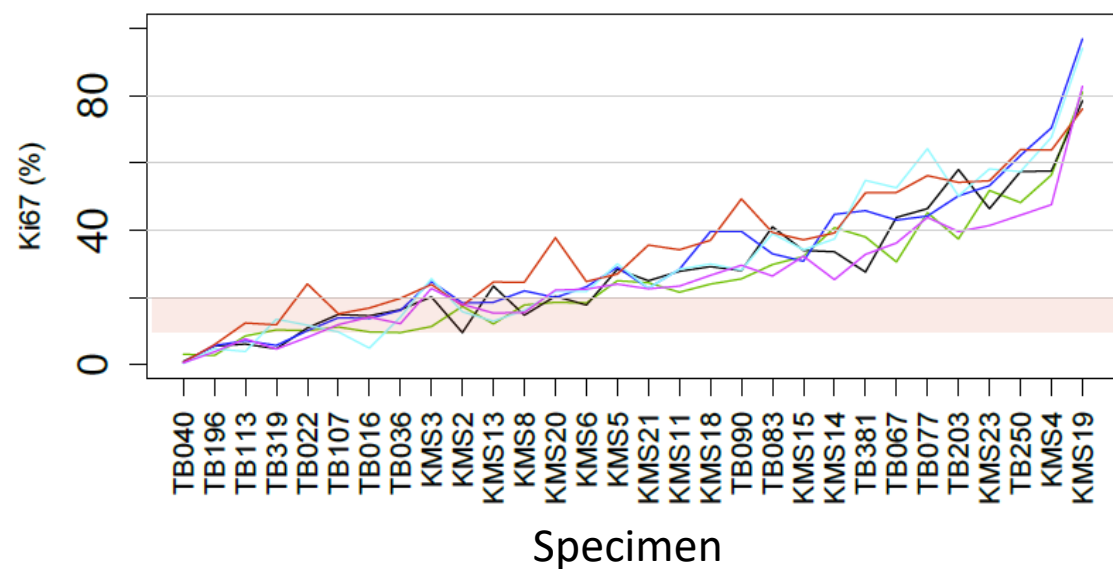


Unweighted global score



Variance component analysis



a**Weighted global score****b****Unweighted global score****c****Hot-spot score**

a. Weighted global score

(median score)		SLIDE SET 1						SLIDE SET 2						SLIDE SET 3					SLIDE SET 4						
		B	F	E	V	G	H	A	P	R	K	T	O	L	J	I	U	M	D	Q	W	N	X		S
Specimen	TB040 (0)	1	0	0	0	1	1	0	1	1	0	0	1	0	0	1	1	0	1	0	0	0	1	0	23:0:0
	TB196 (3)	3	3	6	3	4	6	2	2	5	2	3	4	3	1	4	2	4	2	4	3	4	2	2	23:0:0
	TB113 (6)	6	4	6	6	6	18	4	6	5	6	8	8	9	3	6	6	7	5	4	5	3	6	7	22:1:0
	TB319 (6)	6	6	7	3	10	10	2	5	5	5	6	8	7	4	4	6	6	9	6	7	6	8	4	21:2:0
	TB107 (6)	4	5	5	8	5	10	3	3	9	6	8	6	8	3	8	6	8	7	5	5	4	7	9	22:1:0
	KMS13 (7)	7	5	8	3	8	9	2	2	8	4	11	12	9	4	10	17	11	4	8	7	3	4	6	18:5:0
	TB016 (8)	6	8	8	8	10	15	5	7	8	9	11	13	21	8	13	10	10	10	6	8	7	12	7	13:9:1
	KMS18 (9)	6	8	9	8	10	18	3	12	9	11	16	6	19	6	8	9	10	16	10	5	2	2	5	14:9:0
	KMS3 (9)	7	9	14	9	16	12	7	7	9	7	9	8	12	9	16	12	15	12	6	11	9	11	9	13:10:0
	TB022 (9)	7	8	10	5	10	15	3	10	6	13	15	5	5	8	10	7	9	7	10	6	11	11	14	12:11:0
	KMS2 (10)	13	9	7	14	11	16	5	9	8	14	13	11	20	9	10	10	9	11	6	9	4	11	10	10:13:0
	KMS8 (10)	11	10	12	9	18	22	9	10	10	18	15	8	22	10	11	17	10	14	8	7	9	9	10	7:14:2
	KMS6 (12)	9	13	6	11	14	22	12	11	11	9	16	19	20	10	9	10	15	24	12	11	8	18	15	5:16:2
	TB036 (12)	8	9	12	12	12	20	4	5	13	11	14	13	21	25	15	14	12	12	9	9	8	13	10	7:14:2
	KMS11 (12)	16	10	12	13	10	23	12	14	10	11	14	10	22	7	14	19	16	7	10	12	11	11	14	2:19:2
	KMS20 (14)	12	16	8	12	11	24	10	9	12	13	16	18	24	11	12	16	15	19	36	13	19	14	17	2:18:3
	KMS21 (15)	11	13	16	14	16	15	11	6	11	21	30	19	24	10	15	15	20	22	11	10	15	13	17	1:18:4
	TB090 (16)	14	15	17	14	15	24	20	16	11	17	15	11	16	11	21	10	22	21	14	14	20	20	22	0:18:5
	KMS15 (17)	22	16	19	17	13	29	8	7	12	19	21	19	21	15	17	14	23	19	16	12	13	16	24	2:15:6
	TB381 (18)	17	21	18	13	26	21	16	10	15	14	20	20	29	11	14	20	28	21	14	10	20	11	21	0:16:7
	KMS14 (19)	25	27	17	19	20	24	10	12	15	15	14	28	28	16	24	15	26	20	15	24	19	18	25	0:14:9
	TB083 (20)	21	18	26	23	21	32	25	10	19	19	24	19	29	15	23	17	18	20	9	16	24	28	17	1:11:11
	KMS5 (25)	20	25	28	20	25	28	18	14	13	25	28	18	31	11	28	28	25	21	18	22	18	29	30	0:9:14
	TB203 (26)	26	24	33	24	30	37	18	19	16	26	25	26	45	31	34	37	33	24	23	34	43	18	32	0:4:19
	TB077 (30)	27	34	26	24	47	39	27	18	19	32	33	28	30	15	31	36	23	45	18	34	30	31	30	0:4:19
	TB067 (32)	32	41	33	28	32	36	24	25	25	32	46	24	45	28	31	23	38	33	29	30	29	34	33	0:0:23
	KMS23 (34)	35	28	35	30	47	49	24	25	32	30	38	35	46	34	39	43	45	34	30	34	34	34	40	0:0:23
	TB250 (40)	40	47	46	34	39	47	42	32	37	39	44	41	51	40	50	46	50	35	27	32	40	31	40	0:0:23
	KMS4 (42)	44	42	33	38	47	37	29	28	50	45	28	47	50	31	51	45	53	39	25	48	27	49	33	0:0:23
	KMS19 (68)	81	75	68	67	83	64	88	60	74	70	68	62	64	61	68	81	77	82	62	58	61	62	67	0:0:23

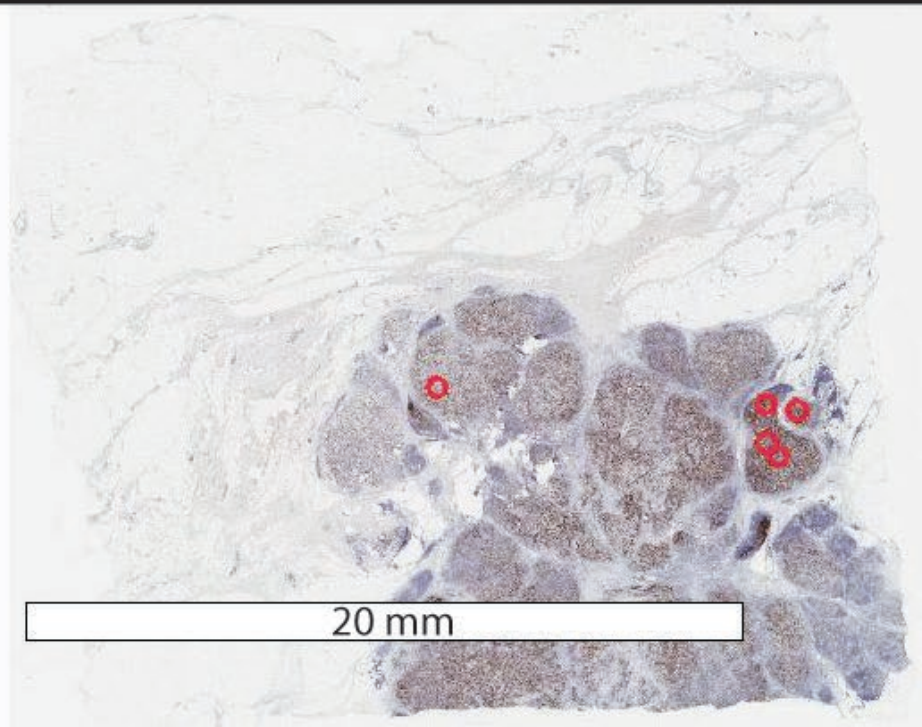
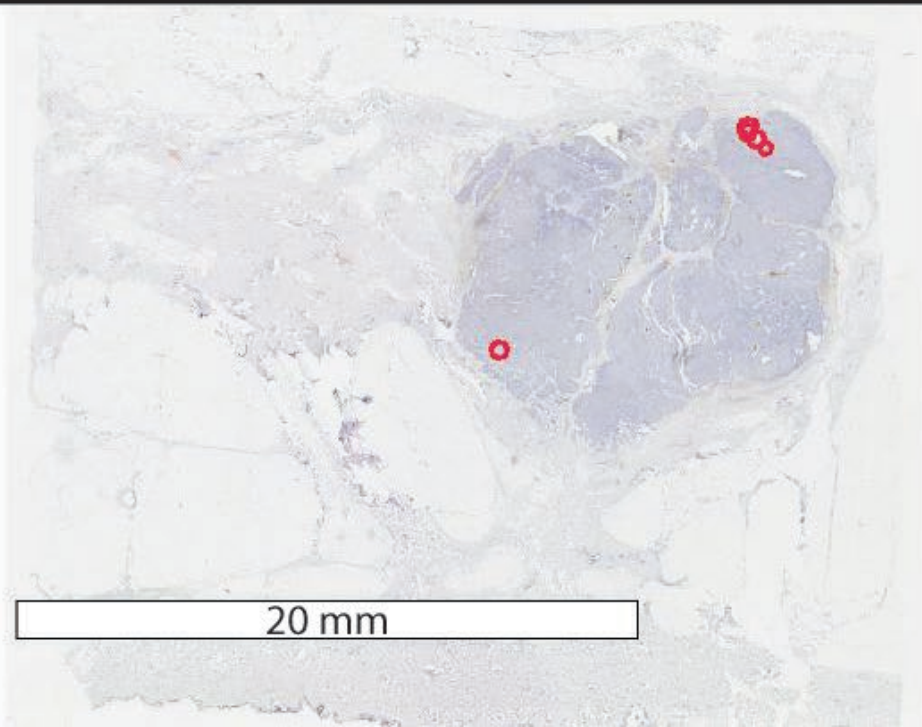
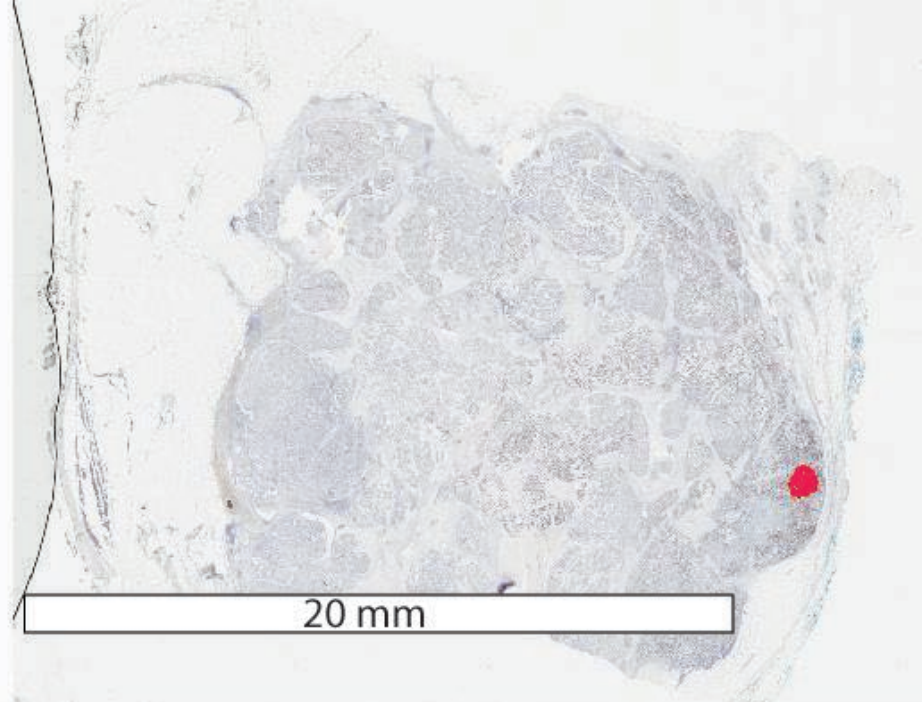
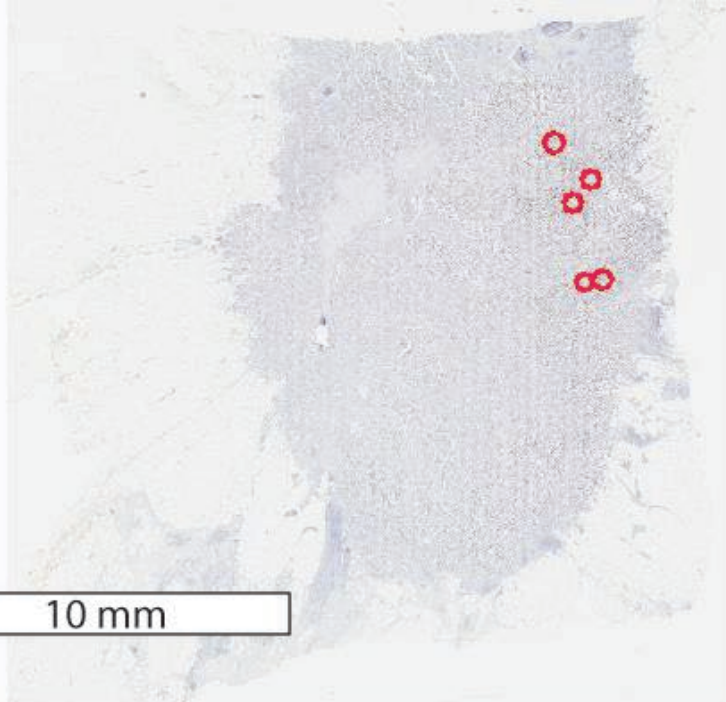
b. Unweighted global score

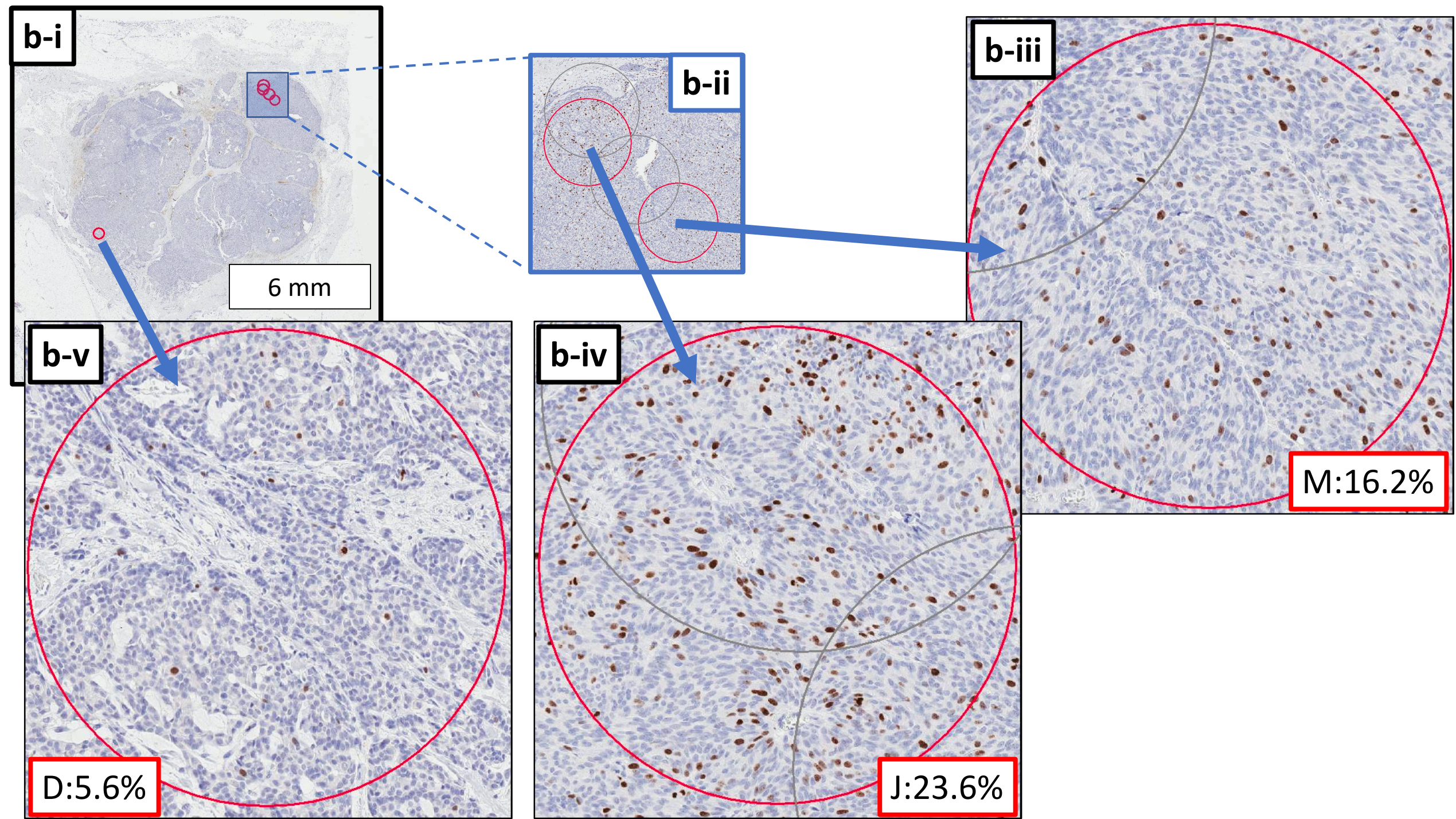
		SLIDE SET 1						SLIDE SET 2						SLIDE SET 3					SLIDE SET 4						
Specimen	(median score)	F	B	V	G	E	H	A	P	R	O	T	K	L	J	I	U	M	D	W	Q	N	X	S	
	TB040 (1)	0	2	0	1	1	1	0	1	1	1	2	0	2	0	2	1	0	1	2	1	0	2	0	23:0:0
	TB196 (4)	6	3	3	4	6	5	2	2	5	5	4	2	2	4	4	2	4	2	3	4	4	5	4	23:0:0
	TB319 (6)	7	6	3	8	8	10	2	4	5	9	7	5	8	4	4	6	7	8	6	6	5	8	4	22:1:0
	TB113 (6)	4	6	6	7	6	9	4	7	5	8	8	6	8	4	6	6	7	5	5	4	3	6	8	23:0:0
	TB107 (6)	6	5	8	6	6	11	3	3	10	6	8	8	8	2	8	6	10	7	5	5	4	7	10	19:4:0
	TB016 (8)	8	5	8	10	9	15	6	7	8	13	11	8	21	7	13	10	10	9	8	6	6	12	8	14:8:1
	KMS13 (8)	8	9	8	11	11	10	4	5	8	13	15	6	14	6	12	15	12	4	8	8	5	4	11	13:10:0
	TB022 (10)	10	8	5	9	10	14	4	10	6	4	14	12	8	10	10	7	11	8	6	10	11	11	14	10:13:0
	KMS2 (10)	9	18	13	12	7	16	5	10	8	11	13	15	21	10	10	9	10	12	10	6	5	11	11	7:15:1
	KMS3 (10)	10	7	10	16	14	11	5	6	9	8	12	8	14	10	15	12	16	11	11	6	9	11	10	8:15:0
	KMS6 (11)	11	10	11	14	6	21	14	9	10	19	17	10	18	11	9	10	16	21	10	12	7	17	18	4:17:2
	TB036 (12)	9	8	12	12	12	19	4	5	12	13	14	10	23	29	15	14	12	12	9	9	8	13	12	7:14:2
	KMS8 (12)	12	12	8	13	16	21	9	10	11	10	12	18	23	12	11	14	12	14	7	8	9	10	14	5:16:2
	KMS11 (12)	14	15	14	13	19	24	12	11	12	10	13	12	27	11	12	20	13	8	12	12	10	11	15	1:20:2
	KMS18 (14)	14	8	14	9	17	20	4	16	15	11	16	15	20	12	12	8	16	16	6	24	4	2	12	7:15:1
	KMS21 (15)	15	10	14	15	16	14	11	6	11	18	26	19	27	12	14	15	19	21	11	13	15	13	20	1:19:3
	KMS20 (15)	16	12	15	11	10	24	12	10	12	17	13	16	28	14	11	16	16	18	13	35	18	14	19	0:20:3
	TB090 (18)	19	15	16	15	24	25	18	16	10	12	15	18	17	16	21	12	20	23	16	18	19	20	28	0:18:5
	KMS15 (18)	13	22	18	18	21	26	8	10	19	19	25	18	22	16	16	14	22	19	15	16	16	17	26	1:15:7
TB381 (18)	21	26	15	25	18	22	19	10	17	15	23	18	32	14	13	18	30	22	12	16	20	16	22	0:14:9	
KMS14 (20)	22	27	18	22	18	20	9	14	15	26	15	18	30	15	24	15	26	20	24	16	23	20	28	1:12:10	
TB083 (21)	19	18	24	21	27	28	22	14	20	16	22	22	32	16	20	14	20	21	18	10	24	25	23	0:11:12	
KMS5 (22)	26	20	21	25	28	28	15	14	13	16	24	22	31	11	27	28	24	18	22	20	18	29	29	0:9:14	
TB077 (27)	30	26	20	44	27	31	19	22	19	28	34	24	32	18	27	36	22	28	32	24	30	28	23	0:4:19	
TB203 (28)	20	29	28	30	37	36	23	20	20	24	27	24	40	27	34	33	33	25	34	24	34	22	31	0:3:20	
TB067 (30)	38	30	26	32	32	28	24	22	25	25	33	30	42	26	30	24	36	30	30	31	30	34	29	0:0:23	
KMS23 (34)	28	37	34	47	36	46	24	25	31	32	38	30	46	32	37	38	44	35	32	27	33	34	43	0:0:23	
KMS4 (36)	31	42	37	47	35	31	32	27	50	40	33	34	44	28	48	38	56	36	49	28	24	51	36	0:0:23	
TB250 (37)	25	34	36	38	44	40	42	31	36	39	36	34	48	37	50	46	45	36	30	22	38	33	37	0:0:23	
KMS19 (66)	72	76	69	72	66	62	88	55	69	64	64	67	64	59	68	81	70	82	59	59	61	62	64	0:0:23	

c. Hot-spot score

		SLIDE SET 1						SLIDE SET 2						SLIDE SET 3					SLIDE SET 4						
Specimen	(median score)	B	V	E	F	G	H	A	P	R	O	L	K	T	U	M	D	J	I	Q	N	W	X	S	
	TB040 (1)	3	1	1	1	0	1	1	1	6	0	2	3	15	0	1	1	1	6	1	0	2	1	1	22:1:0
	TB196 (6)	3	4	6	6	5	6	3	4	7	3	7	6	14	0	7	7	8	11	1	5	7	8	10	20:3:0
	TB113 (8)	9	8	6	7	4	12	9	5	9	7	7	16	10	9	12	12	8	5	5	4	6	14	16	16:7:0
	TB319 (9)	10	5	5	6	14	12	5	3	14	7	12	9	14	7	11	15	8	10	4	4	6	10	15	12:11:0
	TB022 (10)	10	8	11	10	12	24	8	9	15	8	22	20	20	10	15	10	14	6	8	8	9	14	16	8:13:2
	TB107 (12)	11	12	15	14	10	15	5	7	13	7	14	11	19	9	13	11	15	11	5	6	12	14	18	6:17:0
	TB016 (14)	10	14	15	14	5	17	11	9	12	13	24	15	15	10	9	19	20	17	7	11	10	17	14	4:18:1
	TB036 (14)	10	12	16	16	14	20	7	8	15	13	24	19	14	7	14	19	14	13	11	9	11	18	19	4:18:1
	KMS3 (14)	11	23	20	25	26	24	8	12	5	16	29	24	14	9	20	12	20	14	4	9	9	11	22	6:10:7
	KMS2 (16)	17	18	10	18	16	18	9	13	13	11	23	27	24	15	13	17	17	16	9	8	8	15	20	4:16:3
	KMS13 (18)	12	15	23	19	13	25	12	18	21	18	31	26	31	18	16	6	24	21	2	14	13	17	15	2:13:8
	KMS8 (18)	18	16	15	22	16	25	17	15	21	18	20	27	26	18	23	21	13	20	6	12	10	14	19	1:15:7
	KMS20 (22)	19	22	20	20	22	38	17	20	20	23	35	27	29	22	22	27	25	29	21	17	22	23	28	0:7:16
	KMS6 (22)	18	22	18	23	22	25	13	16	12	27	25	25	32	22	26	26	26	14	8	15	18	18	30	1:9:13
	KMS5 (24)	25	24	28	29	30	27	23	15	14	23	33	35	37	28	21	24	22	34	24	17	20	23	33	0:4:19
	KMS21 (24)	24	23	25	23	23	36	14	23	19	23	27	30	31	16	26	23	31	27	14	16	26	25	30	0:5:18
	KMS11 (26)	22	23	28	28	28	34	14	12	29	17	35	33	39	29	21	26	25	26	15	21	30	25	33	0:4:19
	KMS18 (28)	24	27	29	40	30	37	23	27	29	28	28	40	39	7	22	35	35	38	18	21	20	19	33	1:3:19
	TB090 (30)	26	30	28	40	28	49	25	28	30	27	32	39	50	30	26	32	26	36	31	22	28	30	31	0:0:23
	TB083 (31)	30	26	41	33	39	39	22	28	31	33	47	41	46	21	25	33	28	37	27	31	25	31	32	0:0:23
	KMS15 (32)	32	32	34	31	34	37	13	26	30	27	25	39	39	30	26	32	32	27	27	18	32	37	37	0:2:21
	KMS14 (34)	41	25	34	45	37	39	19	22	21	30	34	38	43	30	33	21	35	42	23	21	26	38	38	0:1:22
	TB381 (35)	38	33	28	46	55	51	27	27	31	33	46	41	46	38	35	39	34	39	32	28	26	30	37	0:0:23
TB067 (43)	31	36	44	43	53	51	31	37	49	43	40	46	52	43	49	37	39	54	39	32	47	36	46	0:0:23	
TB077 (45)	45	44	46	44	64	56	25	36	43	31	50	50	49	44	52	46	34	41	43	33	52	49	48	0:0:23	
TB203 (47)	37	40	58	50	50	54	39	42	42	41	47	50	50	45	48	40	48	51	44	47	48	46	52	0:0:23	
KMS23 (48)	52	41	46	53	58	55	42	34	42	36	52	48	55	55	52	41	51	47	46	43	43	48	57	0:0:23	
TB250 (53)	48	44	57	62	57	64	32	43	48	47	51	60	54	56	54	69	53	57	40	48	49	53	49	0:0:23	
KMS4 (56)	56	48	58	70	68	64	38	51	64	45	58	68	54	48	53	53	51	63	55	39	57	61	63	0:0:23	
KMS19 (76)	81	83	78	97	94	76	92	68	66	65	64	79	74	71	81	82	67	84	64	66	81	56	70	0:0:23	

a





Supplemental table 1. Cohort characteristics (from pathology reports and case notes) of the 30 cases.

Age (years)		
	median / mean	61 / 62
	minimum / maximum	40 / 85
Grade		
	grade 1 (%)	5 (17%)
	grade 2 (%)	18 (60%)
	grade 3 (%)	7 (23%)
Tumor size (mm)		
	median / mean	20 / 22
	minimum / maximum	12 / 42
Number of positive nodes		
	0 (%)	21 (70%)
	>0 (%)	9 (30%)
ER (IHC)		
	positive	30 (100%)
	H-score (UK cases only) minimum / maximum	130 / 240
	Percent positive (Japan cases only) minimum / maximum	80 / 100

Supplemental table 2. Summary statistics for weighted global scores (0-100%)¹, ordered according to observer median.

Group 1²							
Observer	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	SD
B	0.9	6.6	12.4	17.8	24.1	80.8	16.4
F	0.1	8.3	12.8	18.3	24.6	74.5	16.2
E	0.3	7.8	12.9	18.1	26.1	67.9	14.6
V	0.5	8.1	13	16.4	22.2	66.8	13.5
G	0.8	9.8	14.2	20.5	25.6	82.8	17.5
H	1	15	22.1	24.1	31	63.9	13.9
Group 2							
Observer	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	SD
A	0.2	3.7	9.8	14.7	19.2	87.8	17.1
P	1.2	6.4	9.9	13.1	15.4	59.6	11.7
R	1.2	8.6	11.1	15.9	15.9	74	15
K	0.2	8.8	13.9	18.1	24	70.5	14.8
T	0.4	11.8	15.4	20.3	26.8	67.5	14.3
O	1	8.4	15.5	18.3	22.9	61.5	13.7
L	0.4	12.9	21.8	24.4	30.1	64.5	15.9
Group 3							
Observer	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	SD
J	0.2	7.4	10.6	14.9	16.1	61.3	13.5
I	1.4	10.2	14.3	19.9	27.1	68.5	15.8
U	1	9.7	14.9	19.7	22.4	80.8	16.9
M	0.5	9.9	15.4	21.3	25.9	76.8	17
D	0.7	9	19.2	19.8	23.3	82.2	16.3
Group 4							
Observer	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	SD
Q	0.4	6.6	10.6	14.9	18.3	61.9	12.5
W	0.3	6.7	11.3	16.5	23.1	57.5	14
N	0.2	5.8	11.9	16.7	22.8	60.8	14.3
X	1	9.5	13	17.8	26.4	61.8	14.1
S	0.2	9	15.7	19	28.5	67	14.5
Mean³							
	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	SD
	13.1	16.4	18.1	18.3	19.9	24.4	2.8

¹All summary statistics are computed using the untransformed Ki67 scores.

²Groups of observers are defined by which collection of slides they received for scoring.

³Mean for each observer computed using untransformed Ki67 scores for 30 slides; summary statistics provided for that collection of observer means.

Supplemental table 3. Summary statistics for unweighted global scores (0-100%)¹, ordered according to observer median.

Group 1²							
Observer	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	SD
F	0.2	9.1	13.5	17.5	21.3	71.8	13.6
B	2	7.5	13.8	18.2	26.2	76.5	15.5
V	0.5	8.5	14.2	17.2	21.1	68.8	13.7
G	0.8	10	14.6	20.2	25.2	71.7	15.8
E	0.8	9	16.6	19.5	26.9	66.5	14.2
H	1.2	13.8	20.9	22.3	27.7	61.8	12.6
Group 2							
Observer	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	SD
A	0.2	4.2	10	14.8	19.2	87.8	17
P	1.2	6.2	10.3	13.4	16	55	10.9
R	1.2	8.6	11.6	16.4	19	69.2	14.2
O	1	10.5	14.2	17.8	22.6	63.5	13
T	2	12.4	15.1	19.8	26	64.5	12.9
K	0.5	8.6	16.9	17.6	22.2	67.2	12.9
L	1.8	15.2	22.9	25	31.7	64.5	14.8
Group 3							
Observer	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	SD
J	0.5	10	12.4	15.8	17.2	59.2	12.3
I	2	10.7	13.8	19.5	26.1	68.5	15.3
U	1	9.4	14.5	18.9	22.7	80.8	16.3
M	0.5	11	16.4	21.3	25.1	70	15.7
D	0.8	8.5	18	19.1	23	82.2	15.5
Group 4							
Observer	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	SD
W	2	7.1	11.5	16.9	23.1	59.2	13.9
Q	0.8	6.6	12.9	15.9	23	58.8	12.1
N	0.2	5	13.1	16.4	23.5	60.8	13.6
X	1.5	9.9	13.4	18.1	24	61.8	14
S	0.5	11.2	18.7	20.2	28.2	63.5	13.4
Mean³							
	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	SD
	13.4	16.7	18.1	18.3	19.7	25	2.6

¹All summary statistics are computed using the untransformed Ki67 scores.

²Groups of observers are defined by which collection of slides they received for scoring.

³Mean for each observer computed using untransformed Ki67 scores for 30 slides; summary statistics provided for that collection of observer means.

Supplemental table 4. Summary statistics for hot-spot scores (0-100%)¹, ordered according to observer median.

Group 1²							
Observer	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	SD
B	2.8	11.3	22.8	25.8	36.1	81	17.9
V	0.6	14.5	23.1	25.3	32.7	82.7	16.7
E	1	14.8	26.3	28.1	39.2	78.4	18.7
F	0.6	16.8	26.6	31.1	43.9	96.8	21.4
G	0.4	13.8	26.9	30.9	47.4	94	22.6
H	0.7	20.8	34.9	34.3	50.6	76	18.8
Group 2							
Observer	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	SD
A	0.6	8.8	15.6	20.1	24.8	92	17.5
P	0.8	9.9	19.1	22	28	68.2	15.7
R	4.6	13.5	20.7	25.7	31.2	65.6	16.5
O	0.4	12.8	23.3	23.9	32.2	64.8	14.9
L	2.5	23	30.2	31.4	44.5	63.6	15.4
K	3.4	20.9	31.4	33.1	40.6	79.4	18
T	10.4	19.6	34.3	34.4	48	73.8	16.4
Group 3							
Observer	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	SD
U	0.4	9.2	21.4	24.9	35.8	70.6	18.1
M	0.8	14.4	22.5	27.5	34.7	80.6	18
D	1.2	15.4	24.8	27.9	36.5	82.2	18.1
J	0.8	15.5	25.4	27.5	34.5	67	15.5
I	5	13.8	26.7	29.9	40.5	83.8	19.3
Group 4							
Observer	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	SD
Q	0.6	6	16.5	21.2	32.1	64.2	17.6
N	0.4	8.8	16.9	20.9	30.1	65.6	15.6
W	1.8	10.4	21	25.1	31.1	80.8	18.9
X	1.4	14.3	23	26.7	36.4	60.6	15.6
S	0.8	18	30.4	31	37.8	69.6	16.6
Mean³							
	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	SD
	20.1	25	27.5	27.3	31	34.4	4.2

¹All summary statistics are computed using the untransformed Ki67 scores.

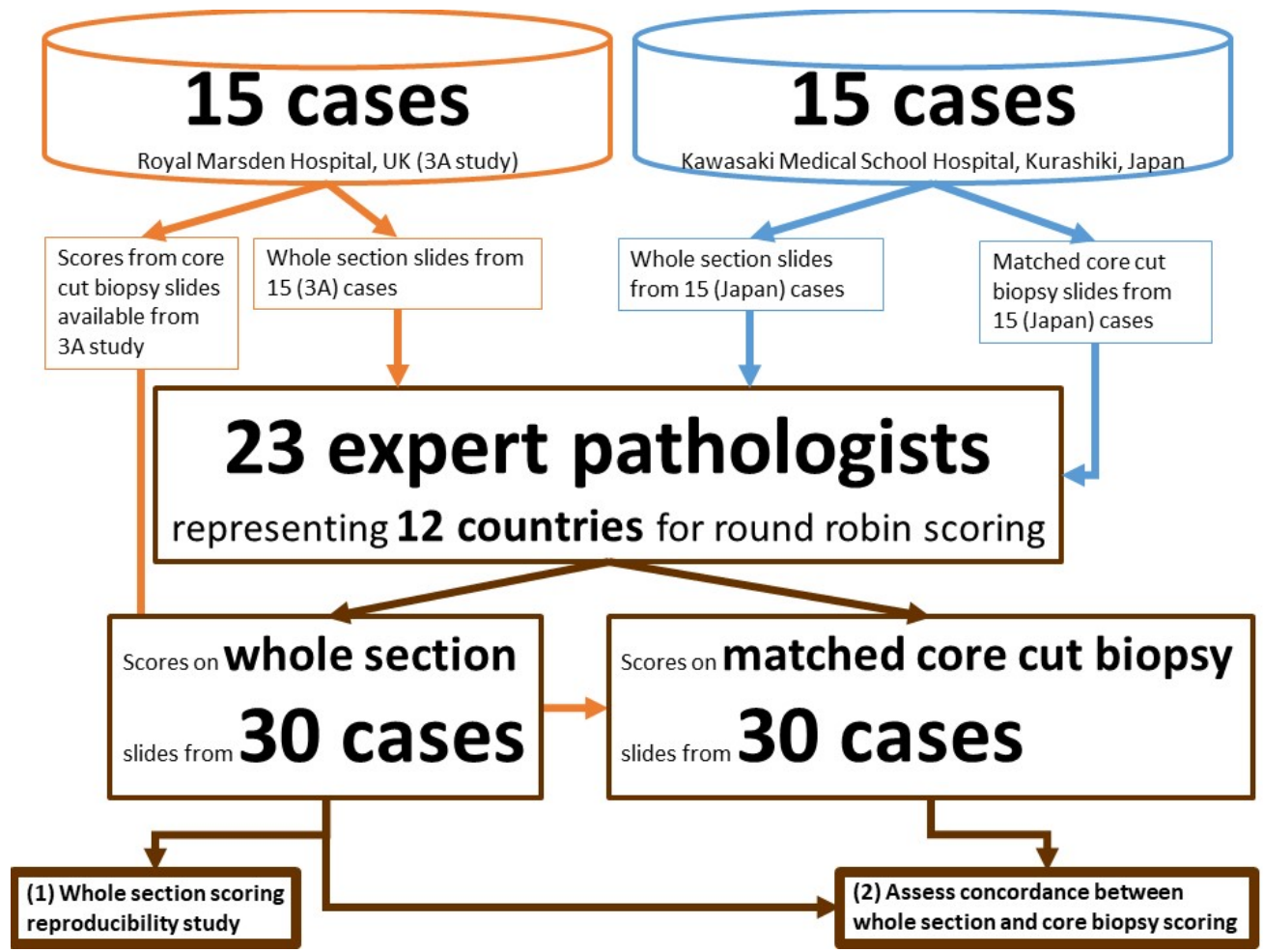
²Groups of observers are defined by which collection of slides they received for scoring.

³Mean for each observer computed using untransformed Ki67 scores for 30 slides; summary statistics provided for that collection of observer means.

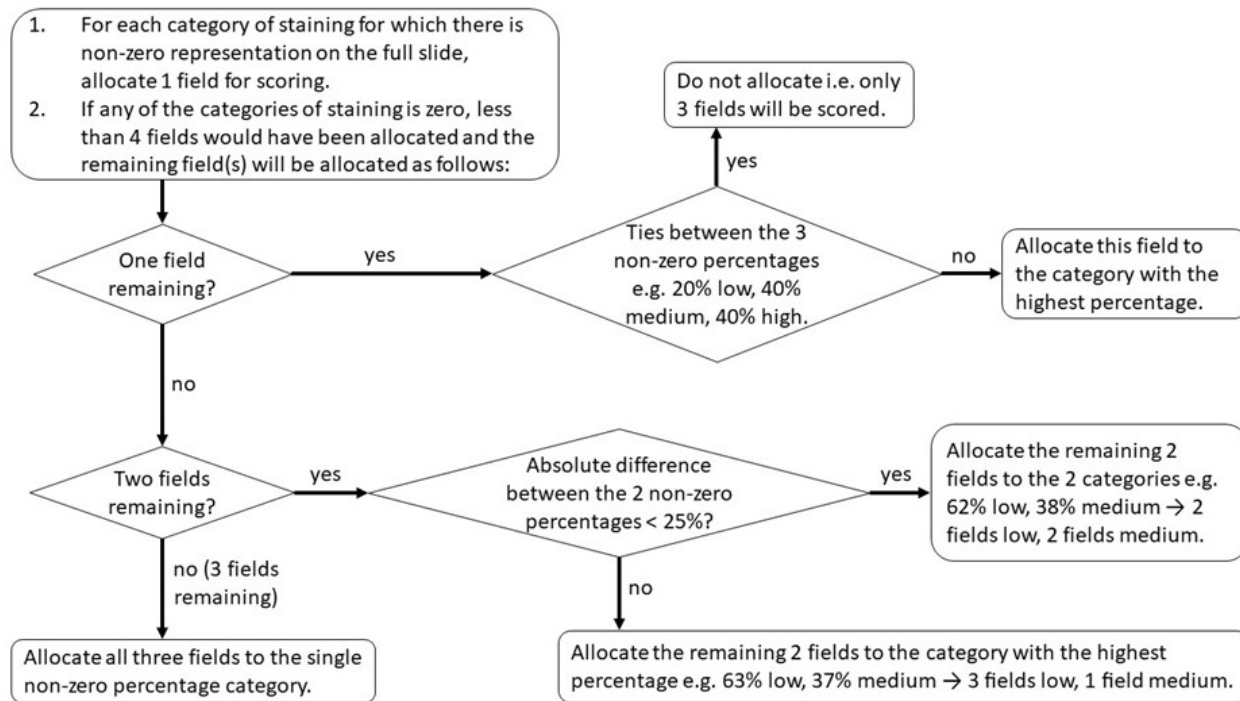
Supplemental table 5. Variance components estimates¹ and corresponding credible intervals.

Component	Scoring method	Variance estimate (95% credible interval)
Biological	Weighted global	1.82 (0.97 – 2.83)
	Unweighted global	1.49 (0.81 – 2.33)
	Hot-spot	1.52 (0.79 – 2.39)
Residual	Weighted global	0.18 (0.17 – 0.21)
	Unweighted global	0.17 (0.15 – 0.18)
	Hot-spot	0.20 (0.17 – 0.22)
Observer	Weighted global	0.08 (0.04 – 0.14)
	Unweighted global	0.06 (0.03 – 0.11)
	Hot-spot	0.10 (0.04 – 0.17)
Section	Weighted global	0.0003 (0 – 0.0014)
	Unweighted global	0.0001 (0 – 0.0004)
	Hot-spot	0.0001 (0 – 0.0002)

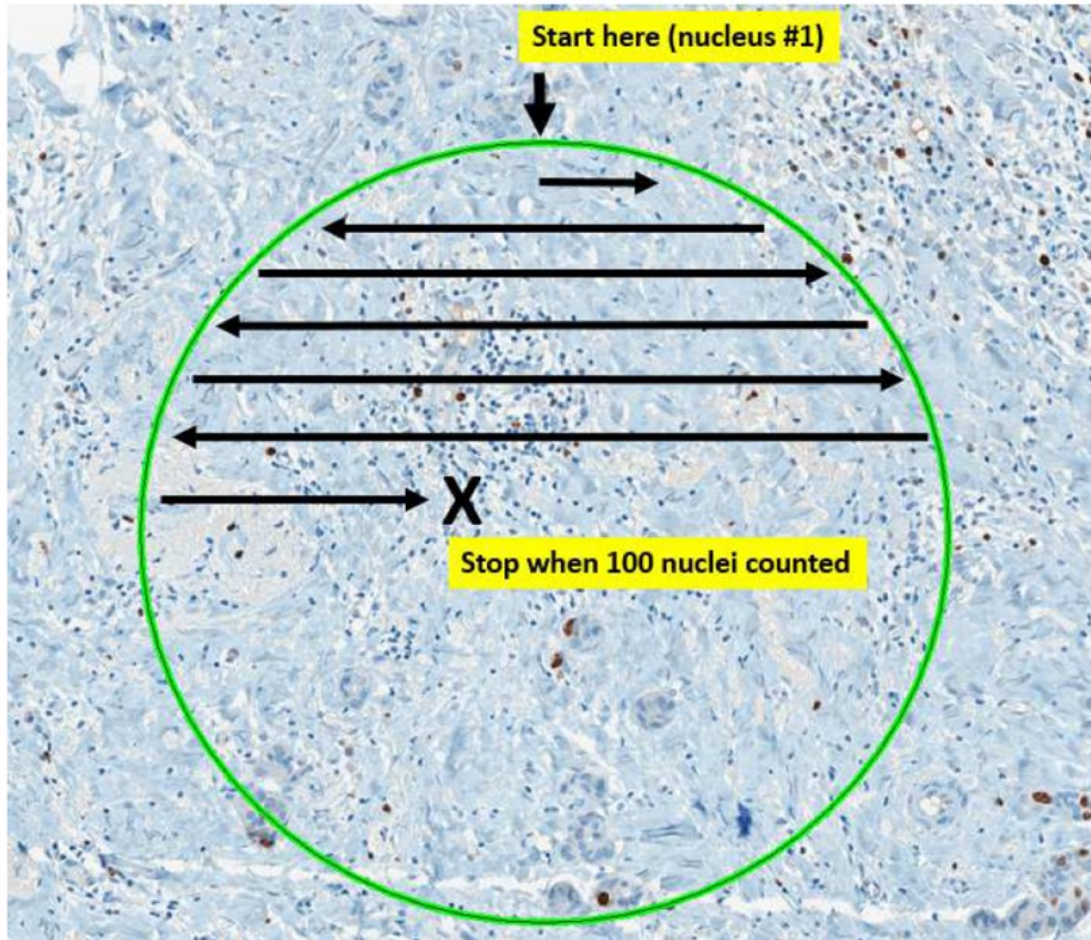
¹Variance component estimates were computed using log₂-transformed Ki67 scores.



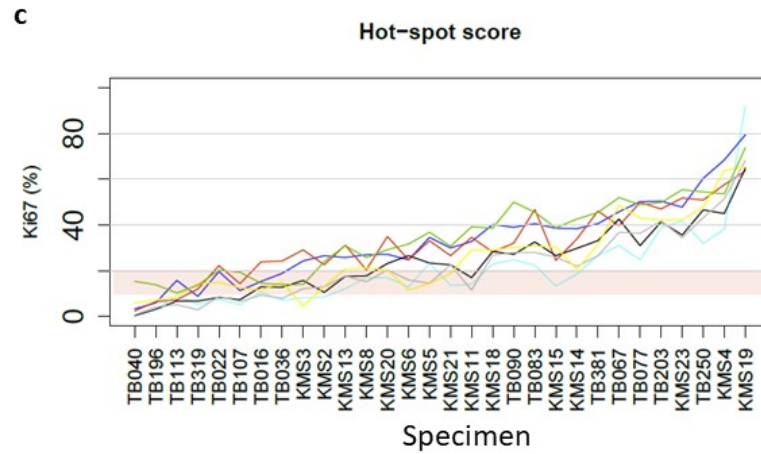
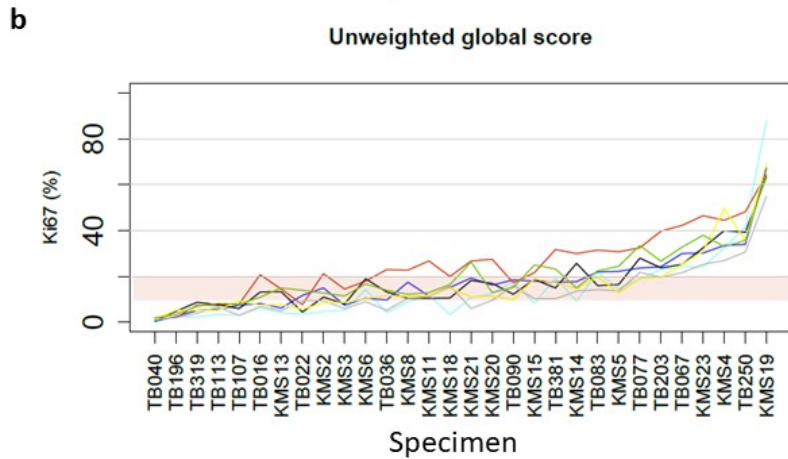
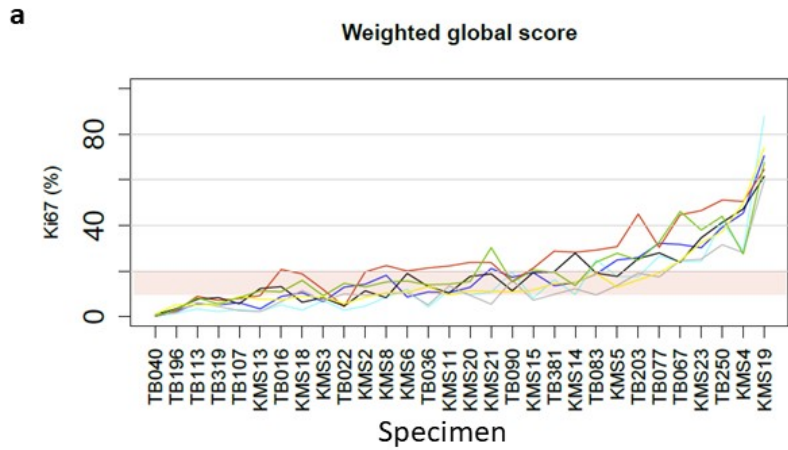
Supplemental Figure 1.
Specimen selection and study design flowchart.



Supplemental Figure 2. Scoring field allocation algorithm.

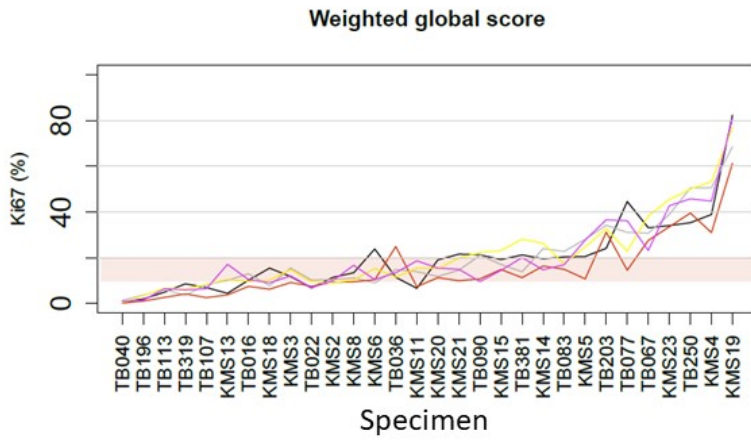


Supplemental Figure 3.
Typewriter nuclei counting
pattern. The green circle indicates
the selected scoring field.

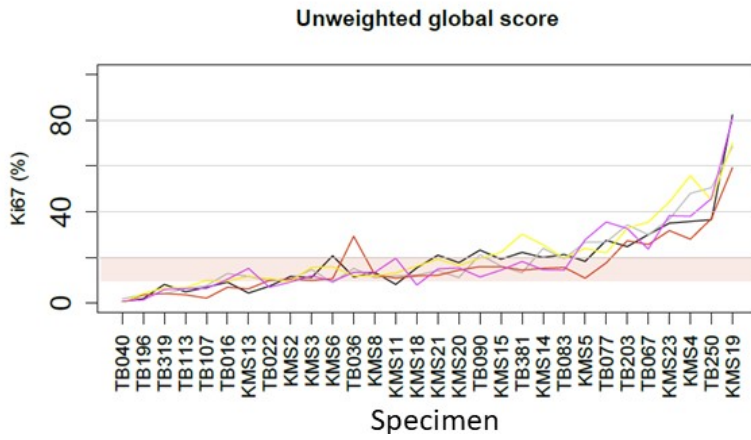


Supplemental Figure 4. Variability in Ki67 scores (a-c correspond to slide set 2, d-f correspond to slide set 3 and g-i correspond to slide set 4). Each line represents Ki67 scores from one observer. Shaded region indicates Ki67 scores between 10-20%.

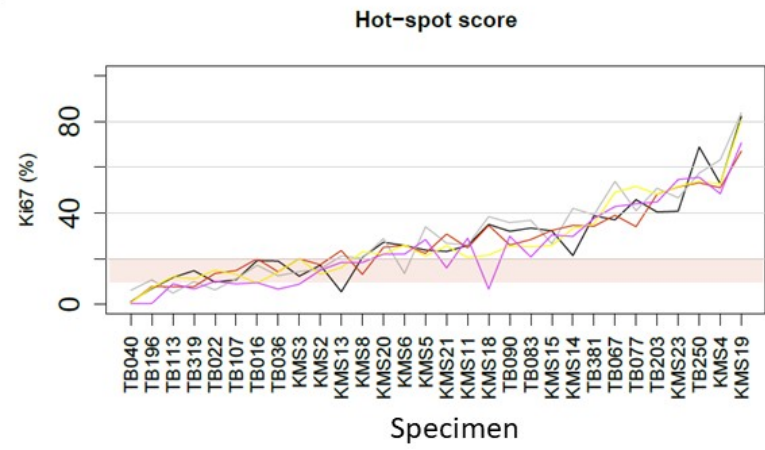
d



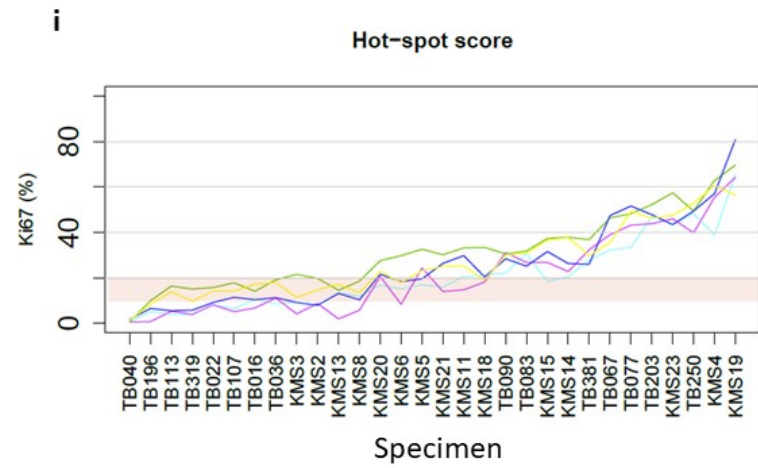
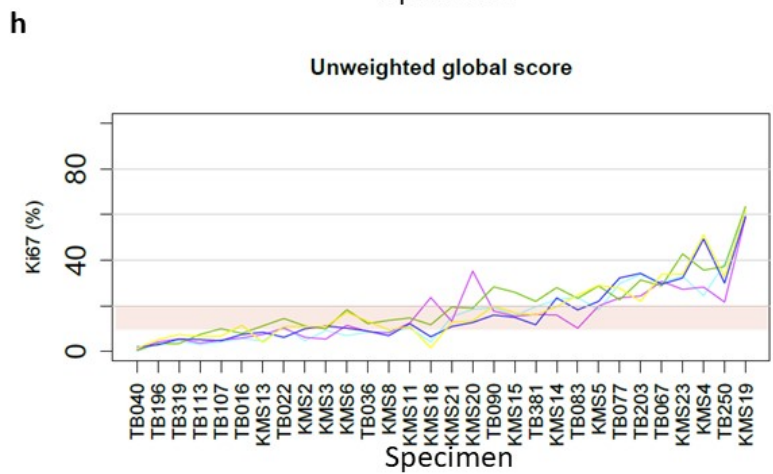
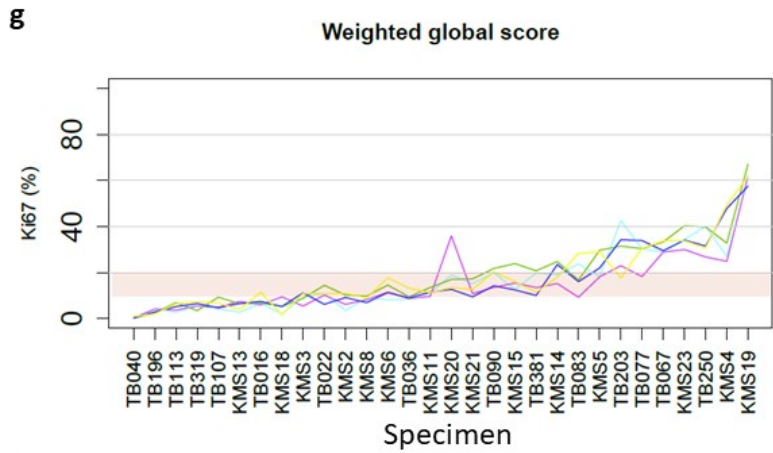
e



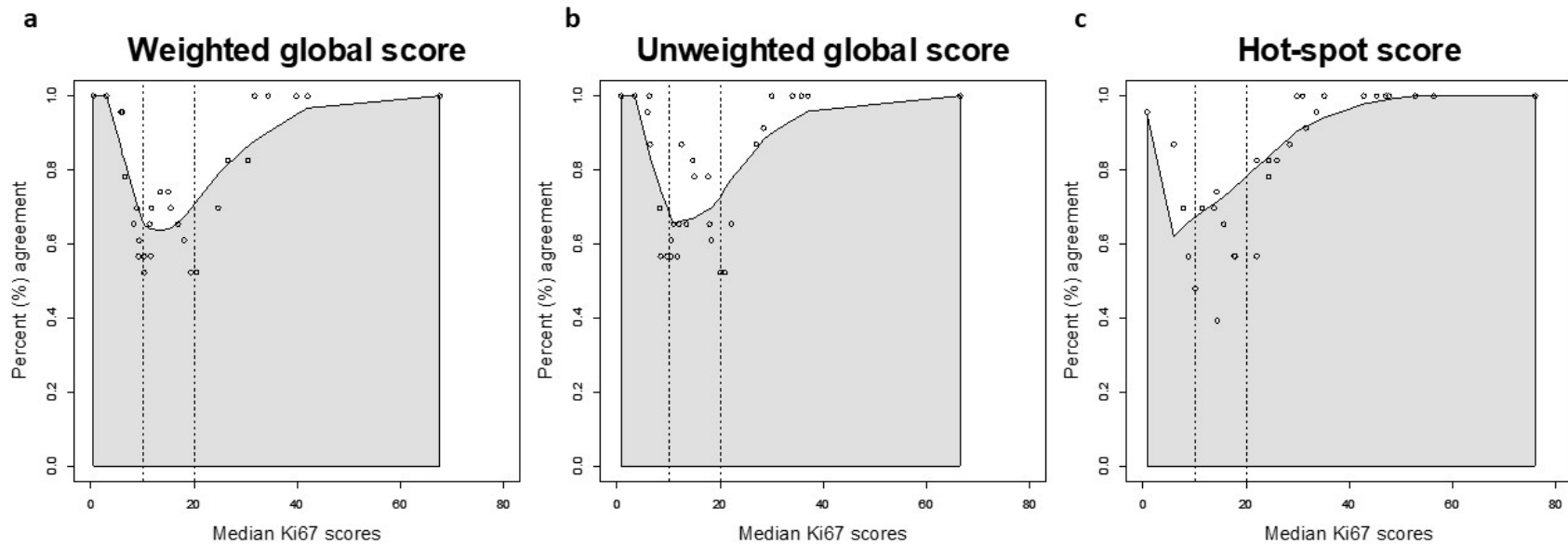
f



Supplemental Figure 4 continued.



Supplemental Figure 4 continued.



Supplemental Figure 5. Percent agreement on categories (<10%, 10–20%, >20%) over Ki67 scores (a, weighted global; b, unweighted global; c, hot-spot). For each of the 30 cases, percentage of observers giving a score that falls into the same category as the median score is plotted against the median score. The dotted lines indicate 10 and 20%. Locally-weighted polynomial regression (performed using the R function “lowess” from the stats package^{1,2}; function was applied to log2-transformed data with default parameters) has been applied to highlight the general trend. Empty circles represent actual data points.

REFERENCES

1. Cleveland WS. LOWESS: A program for smoothing scatterplots by robust locally weighted regression. *The American Statistician*. 1981;35:54.
2. R Core Team. R: A language and environment for statistical computing. R foundation for statistical computing, Vienna, Austria. . 2015.