

Analytical validation of a standardised scoring protocol for Ki67 immunohistochemistry on breast cancer excision whole sections: an international multicentre collaboration

Samuel C Y Leung,¹  Torsten O Nielsen,¹ Lila A Zabaglo,² Indu Arun,³ Sunil S Badve,⁴ Anita L Bane,⁵ John M S Bartlett,^{6,7} Signe Borgquist,⁸ Martin C Chang,⁹  Andrew Dodson,¹⁰ Anna Ehinger,¹¹  Susan Fineberg,¹² Cornelia M Focke,¹³ Dongxia Gao,¹ Allen M Gown,¹⁴ Carolina Gutierrez,¹⁵ Judith C Hugh,¹⁶ Zuzana Kos,¹⁷ Anne-Vibeke Lænkholm,¹⁸ Mauro G. Mastropasqua,¹⁹ Takuya Moriya,²⁰ Sharon Nofech-Mozes,²¹ C Kent Osborne,¹⁵ Frédérique M Penault-Llorca,²² Tammy Piper,⁷ Takashi Sakatani,²³ Roberto Salgado,^{24,25} Jane Starczynski,²⁶ Tomoharu Sugie,²⁷ Bert van der Veegt,²⁸  Giuseppe Viale,^{19,29} Daniel F Hayes,³⁰ Lisa M McShane,³¹ 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)

¹University of British Columbia, Vancouver, BC, Canada, ²The Institute of Cancer Research, London, UK, ³Tata Medical Center, Kolkata, India, ⁴Indiana University Simon Cancer Center, Indianapolis, IN, USA, ⁵Juravinski Hospital and Cancer Centre, McMaster University, Hamilton, ⁶Ontario Institute for Cancer Research, Toronto, ON, Canada, ⁷Edinburgh Cancer Research Centre, Western General Hospital, Edinburgh, UK, ⁸Division of Oncology and Pathology, Department of Clinical Science, Lund University, Lund, Sweden, ⁹Department of Pathology and Laboratory Medicine, University of Vermont Medical Center, Burlington, VT, USA, ¹⁰Ralph Lauren Centre for Breast Cancer Research, The Royal Marsden Hospital, London, UK, ¹¹Department of Clinical Genetics and Pathology, Skane University Hospital, Lund University, Lund, Sweden, ¹²Montefiore Medical Center and the Albert Einstein College of Medicine, Bronx, NY, USA, ¹³Dietrich-Bonhoeffer Medical Center, Neubrandenburg, Germany, ¹⁴PhenoPath Laboratories, Seattle, WA, ¹⁵Lester and Sue Smith Breast Center and Dan L. Duncan Comprehensive Cancer Center, Baylor College of Medicine, Houston, TX, USA, ¹⁶University of Alberta, Edmonton, AB, ¹⁷University of Ottawa and The Ottawa Hospital, Ottawa, ON, Canada, ¹⁸Department of Surgical Pathology, Zealand University Hospital, Slagelse, Denmark, ¹⁹European Institute of Oncology, Milan, Italy, ²⁰Kawasaki Medical School, Kurashiki, Japan, ²¹University of Toronto Sunnybrook Health Sciences Centre, Toronto, ON, Canada, ²²Centre Jean Perrin and Université d'Auvergne, Clermont-Ferrand, France, ²³Nippon Medical School, Bunkyo-ku, 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, UK, ²⁷Kansai Medical University, Hirakata, Japan, ²⁸University Medical Center Groningen, Groningen, the Netherlands, ²⁹University of Milan, Milan, Italy, ³⁰University of Michigan Rogel Cancer Center, Ann Arbor, MI, and ³¹National Cancer Institute, Bethesda, MD, USA

Date of submission 18 December 2018

Accepted for publication 19 April 2019

Published online Article Accepted 24 April 2019

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, Mastropasqua M G, Moriya T, Nofech-Mozes S, Osborne C K, Penault-Llorca F M, Piper T, Sakatani T, Salgado R, Starczynski J, Sugie T, van der Veegt B, Viale G, Hayes D F, McShane L M & Dowsett M

(2019) *Histopathology* 75, 225–235. <https://doi.org/10.1111/his.13880>

Address for correspondence: S C Y Leung, University of British Columbia, Room 509, 2660 Oak Street, Jack Bell Research Center, Vancouver, BC V6H 3Z6, Canada. e-mail: sam.leung@vch.ca

Analytical validation of a standardised scoring protocol for Ki67 immunohistochemistry on breast cancer excision whole sections: an international multicentre collaboration

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 programme to determine whether Ki67 measurement can be analytically validated and standardised among 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: (i) global: four fields of 100 tumour cells each were selected to reflect observed heterogeneity in nuclear staining; (ii) 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 [confidence interval (CI) = 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 min for global and hot-spot methods, respectively.

Conclusions: The global scoring method demonstrates adequate reproducibility to warrant next steps towards 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: analytical validity, immunohistochemistry, interobserver reproducibility, interobserver variability, Ki67, pathology, scoring protocol

Introduction

The nuclear antigen recognised by the Ki67 antibody is expressed in proliferating cells but absent in resting cells.¹ Since its discovery in 1983 by Gerdes *et al.*,¹ Ki67 assessed by immunostaining has been studied extensively as a prognostic^{2–11} and predictive^{4,6,9,12,13} marker, predominantly in hormone receptor-positive breast cancer, but also in other tumours.^{14–18} For example, presurgical Ki67 has been shown to be a marker for recurrence-free survival¹⁹ and, in the neoadjuvant setting, a marker for endocrine-resistant tumour that may require more aggressive treatment.²⁰ Excellent intra-observer reproducibility under controlled pre-analytical 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 management.^{22–24} However, poor interobserver 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 programme to determine whether Ki67 scoring can be standardised and analytically validated throughout

laboratories.^{9,21,29,30} In Phase I, as assessed by the intraclass correlation coefficient (ICC) estimate of interobserver reproducibility, differences in pathologists' visual interpretation were the main source of variability (ICC = 0.71, 95% credible interval (CI) = 0.47–0.78).²¹ Greater concordance was achieved in Phase II, at least on tissue microarrays, when pathologists were trained to calibrate and standardise 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 excision specimens, which require general assessment of the entire sample and selection of areas for formal counting. Therefore, in Phase IIIA, we assessed whether acceptable performance could be achieved on core-cut biopsies using a standardised 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 interobserver 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

The current study represents the final Phase (IIIB) of the visual scoring analytical validity programme,

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 will 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 oestrogen receptor (ER)-positive breast cancer cases were selected: 15 from the Phase IIIA study³⁰ and 15 from Kawasaki Medical School Hospital, Kurashiki, Japan (Figure S1). Case selection was irrespective of patients' age at diagnosis, tumour grade, size or nodal status. The clinicopathological characteristics of these 30 cases are shown in Table S1. All blocks were sectioned and stained in the Royal Marsden Hospital Histopathology Department using monoclonal antibody MIB1 at dilution 1:50 (Dako UK, Ely, UK) using an automated staining system (Ventana Medical Systems, Tucson, AZ, USA) according to criteria established by the IKWG.⁹ Sections from the same block were stained in a single immunohistochemistry run, except for four cases where the staining was performed 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 IIIA study, were invited to participate.

Six adjacent sections from each of the 30 excision blocks were centrally stained: the first with haematoxylin and eosin (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 would have made the study impractical if all pathologists reviewed the same physical glass slides, 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 standardised regions for nuclei counting, through the publicly available proficiency training module (<http://www.gpec.ubc.ca/calibrator>) that was initially used in the Phase II study.²⁹ The detailed scoring protocol is found in Data S1. 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 were the same as those employed in the Phase IIIA study³⁰: (i) 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; (ii) an unweighted global assessment; and (iii) 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 laboratory,³¹ these two global methods require the pathologist to first assess staining heterogeneity by estimating the percentages of the invasive tumour component of the slide exhibiting relatively high, medium, low or negligible Ki67 staining frequencies. Based on these estimates, an algorithm (Figure S2) dictates the required number of fields to select and score for each Ki67 staining frequency (irrespective of staining intensity, totalling up to four fields). This algorithm was designed such that the four (or fewer) 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 tumour nuclei within each field are counted using a 'typewriter' pattern (Figure S3), similar to how

a tissue microarray core was scored in the Phase II study.²⁹

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 tumour nuclei in a 'type-writer' pattern.

STATISTICAL ANALYSES

Prespecified criterion for success

Prior to data collection it was hypothesised 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 II 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 IIIA study.³⁰ 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 tumour nuclei counted divided by the total number of tumour nuclei counted. The weighted global score was derived with tumour 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 IIIA study.³⁰ Briefly, variance component analyses were performed to quantify the contributions from the following sources of variability: scoring pathologist (observer), patient tumour (biological variation – each excision block represents a unique patient) and section of the excision block. Similar to the Phase IIIA 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 generalised linear mixed models.

All data analyses were performed using R version 3.3.2.³³ Sources of variation in log₂-transformed Ki67 scores were analysed using random effects models as implemented in the R packages lme4 and MCMCglmm. Data were visualised using heat maps, box-plots 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 prespecified 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)

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)

ICC, Intraclass correlation coefficient; CI, Confidence interval.

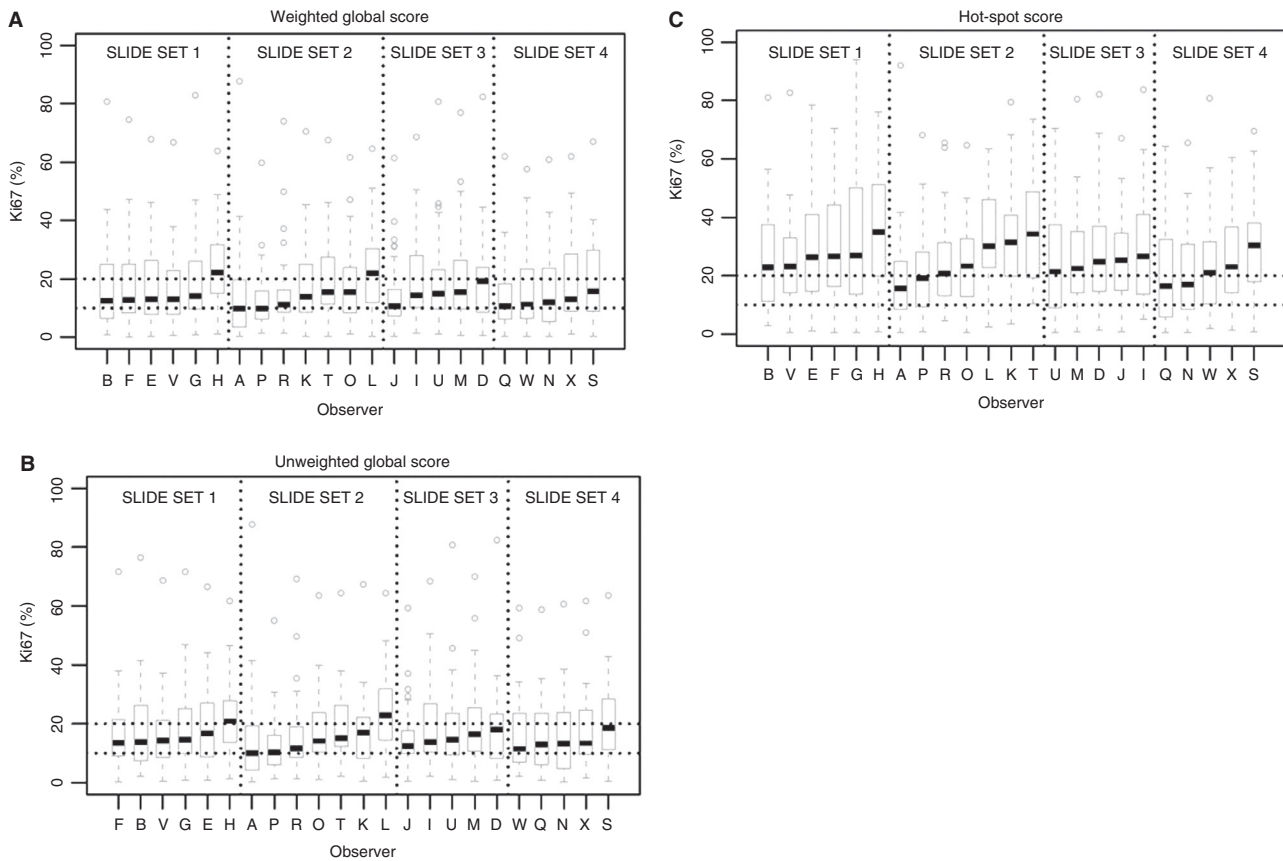


Figure 1. Ki67 scores of all 23 observers (by slide set). Observers are ordered (within each group) by the median scores. The bottom/top of the box in each box plot represent the first (Q1)/third (Q3) quartiles, the bold line inside the box represents the median and the two bars outside the box represent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

respectively, supporting that differences between serial sections were minimal. Figure 1 displays the side-by-side box-plots of Ki67 scores among pathologists (hereafter referred to as 'observers') by group. Summary statistics for the Ki67 scores among the 23 observers are given in Tables S2–S4.

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

In a context where pre-analytical 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 tumour (Figure 2, Table S5).

INTEROBSERVER 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 sets 2–4 are displayed in Figure S4. 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 throughout 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 Figure S5. It

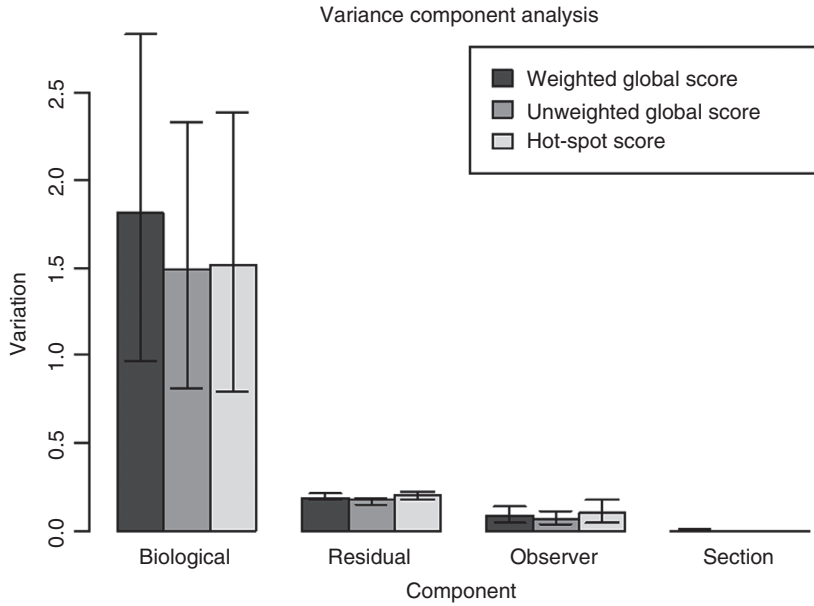


Figure 2. Variance component analysis. Variation due to different components are presented in a bar plot to show the relative magnitude of differences between them. Numerical values of the variance components estimates and the corresponding credible intervals are shown in Table S5.

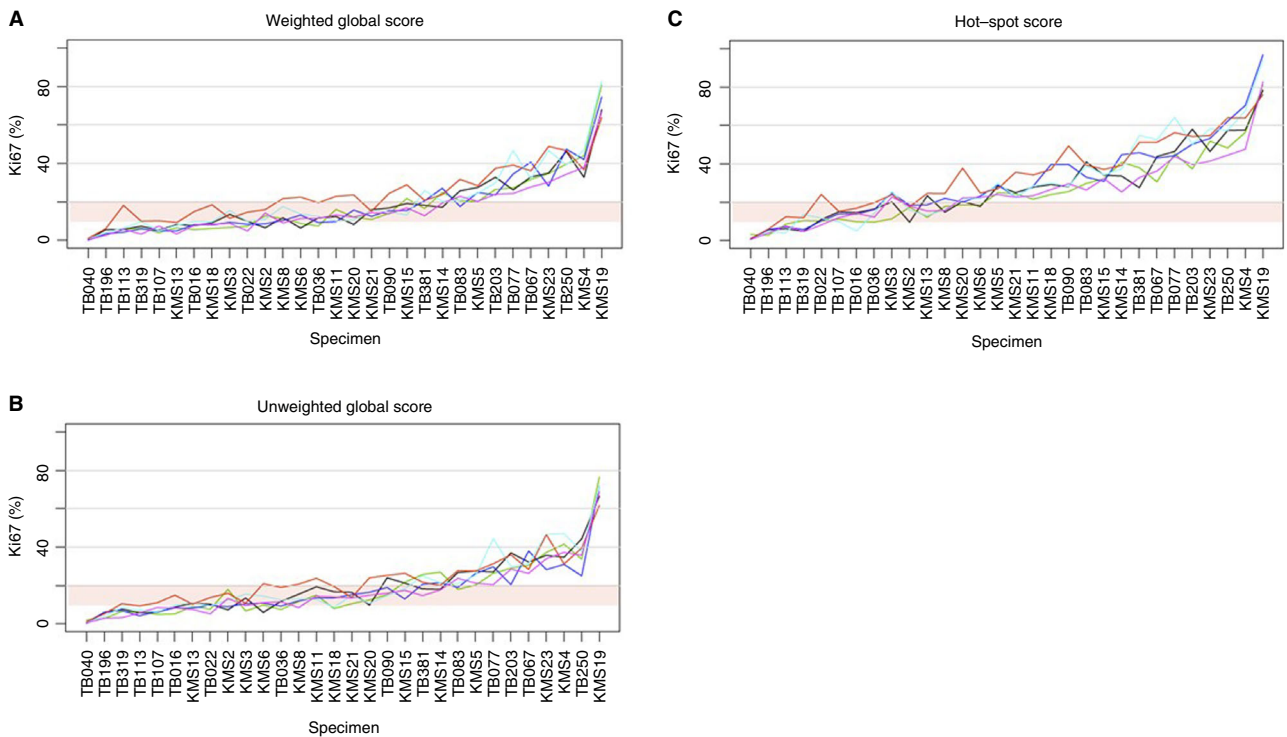


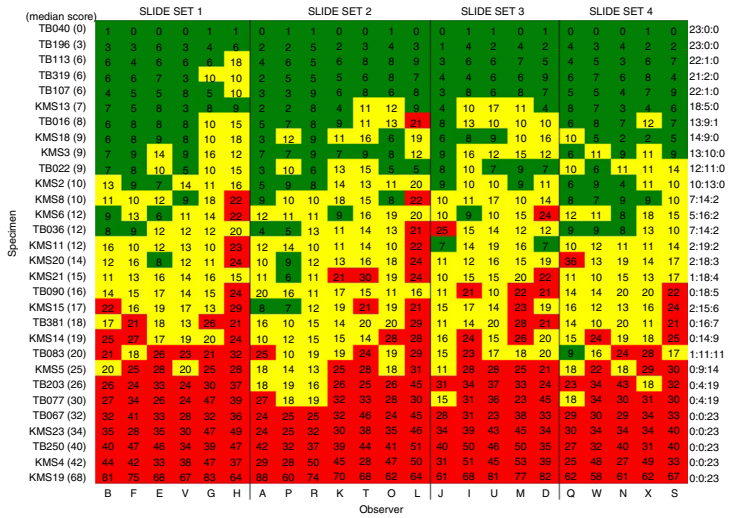
Figure 3. Variability in Ki67 scores (slide set 1 only). Each line represents Ki67 scores from one observer. Shaded region indicates Ki67 scores between 10% and 20%. Scores on slide sets 2–4 are shown in Figure S4.

shows excellent to perfect concordance on cases with scores that are either perfect much lower or higher than the intermediate range (10–20%).

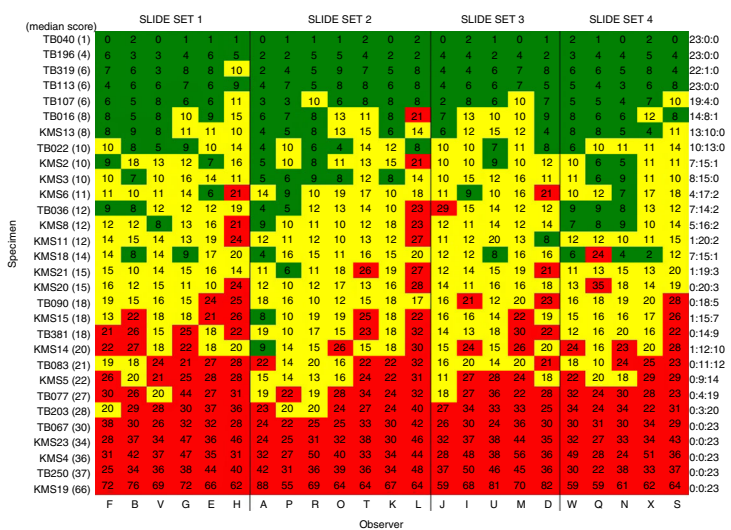
Based on visual inspection of captured images, locations of the hot-spot selections tended to cluster in the

same region among 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>).

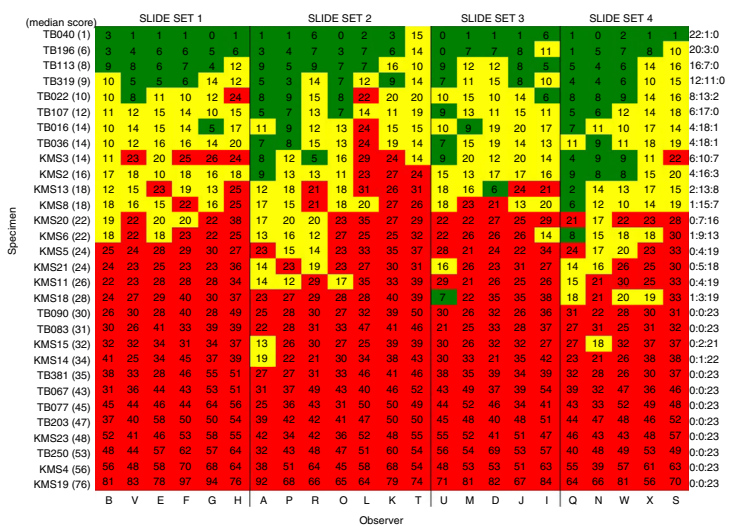
A Weighted global score



B Unweighted global score



C Hot-spot score



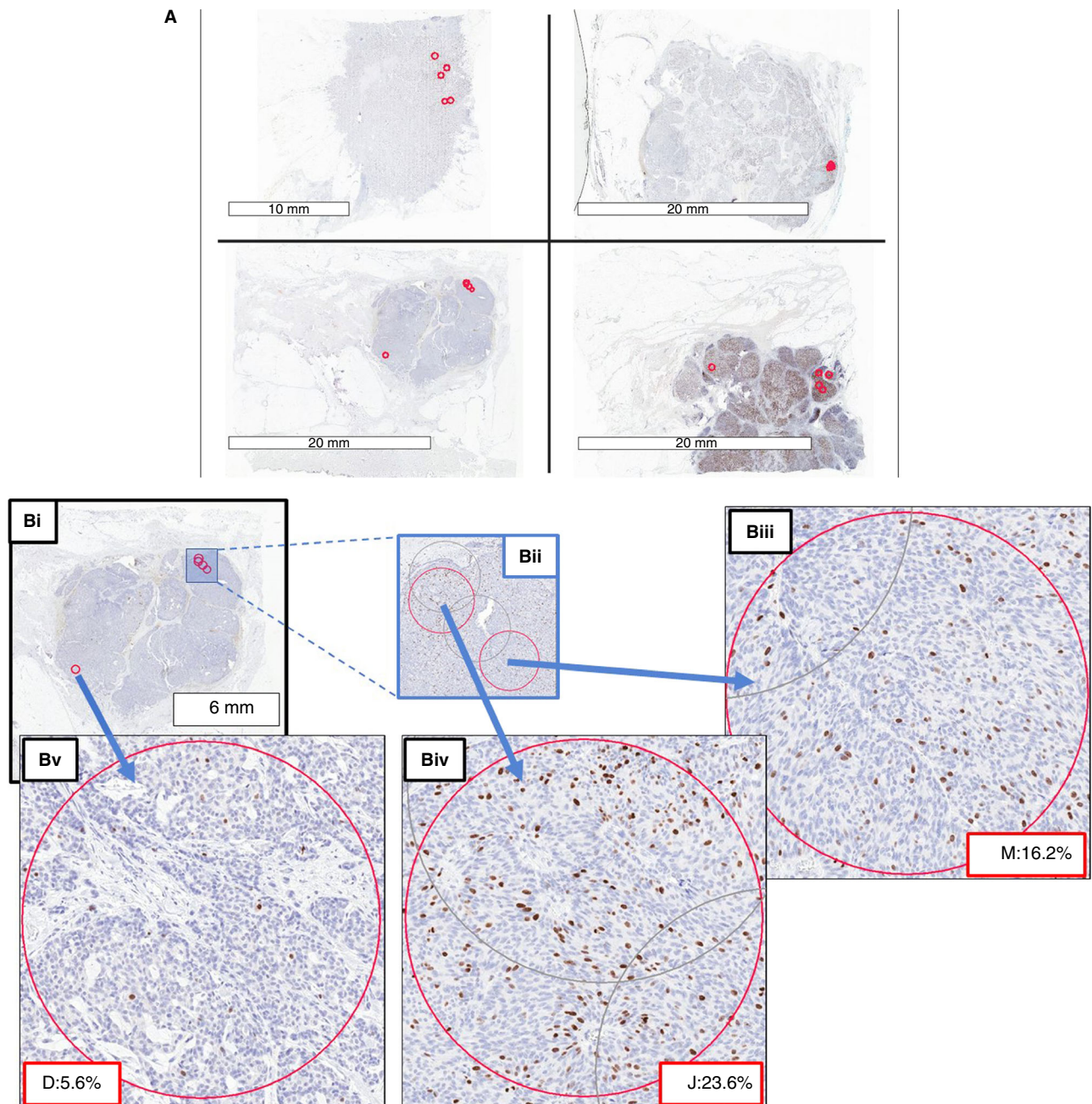


Figure 5. Hot-spot field selection by different observers on the same excision whole section slide. **A,** Selections (indicated by red circles) on some example excision whole section slides. **B,** An example of a single excision whole section slide (median score: 18%) with zoomed-in fields. Each observer 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. B-iii and B-iv represent segments of the same area chosen by two different observers to read Ki67. Figure B-v represents the 'outlier' field selected by only one observer as the hot-spot.

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 pre-

analytical and analytical aspects of the Ki67 immunohistochemistry assay,⁹ and given a set of clearly defined training exercise and scoring instructions, for pathologists to achieve high interobserver 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 a 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 biopsies.³⁰

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 (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 observers 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 approximately 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 min, 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 prespecified (i.e. 400), to determine whether a case has unweighted global score $\geq 10\%$ the pathologist can stop counting if the first field they scored is $\geq 40\%$. For cases with a very low Ki67 score, one would probably still need to count all 400 nuclei.

The proposed scoring protocols do not make any recommendation concerning the required minimum tumour 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 tumour 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.

An external quality assessment programme (e.g. NordiQC³⁴), involving comparison of laboratory scores with reference scores in periodic assessment challenges, will probably improve interobserver reproducibility further. Recent studies suggest that an even higher level of concordance can be achieved with automated image analysis.^{35–38} The IKWG is actively conducting studies in this area to assess how artificial intelligence may help to standardise Ki67 assessment.^{35,38} Also, concordance between Ki67 scores on core biopsies and excision specimens is currently being investigated.

In conclusion, this study demonstrates that an adequately high level of interobserver 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. Pre-analytical and analytical aspects of the immunohistochemistry assay, areas that still need standardisation before the clinical utility of this marker can be proved, will probably 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. 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. J.C.H. 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.

Conflict of interest

S.S.D. has participated in Scientific Advisory Boards/Speaker for Genomic Health Inc., Dako/Agilent, Roche Diagnostics, Targos GmbH, Athenax, Konica-Minolta and received compensation. S.S.D. has received research funding or in-kind support from Dako/Agilent, and has intellectual property right/ownership interests with IU. He is also associated with two startup companies (SYSGenomics and YeS-Genomics). J.M.S.B. has consulted for BioNTech GmbH, Biotheranostics Inc, RNA Diagnostics, and received compensation. He has participated in Scientific Advisory Boards for Biotheranostics and RNA Diagnostics and received compensation and has received research funding or in-kind support from Nanostring, Biotheranostics Inc, BioNTech GmbH. He has intellectual property right/ownership interests with OICR/FACIT. S.B. has participated in educational talks/covered scientific conferences by Roche and Novartis. M.D. 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). A.E. has participated in educational talks organised by Roche but without economical compensation. S.F. participated in a scientific advisory board for Genomic Health and has received monetary compensation (not for salary). D.F.H. 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 OncImmune LLC and InBiomotion, and he serves as a paid advisor for Cepheid, Freenome, CellWorks, Agendia and CVS Caremark. A.-V.L. 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). T.O.N. has consulted for Nanostring and received compensation. He has intellectual property rights/ownership interests from Bioclassifier LLC. C.K.O. has consulted for Astra Zeneca, Genentech and NanoString and received compensation. F.M.P.-L. has participated in Scientific Advisory Boards for Nanostring, Myriad, Genomic Health, Agendia, AstraZeneca, Roche, Sanofi, Novartis, Pfizer, BionTech and received

compensation. He has received research funding or in-kind support from Nanostring, AstraZeneca, Roche, BionTech. B.V.d.V. has consulted for Philips and received compensation. All other authors declare no conflict of interest.

References

- 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; **31**: 13–20.
- 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; **132**: 895–915.
- de Azambuja E, Cardoso F, de Castro G 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; **96**: 1504–1513.
- 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; **24**(Suppl. 2): S67–S72.
- 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; **139**: 539–552.
- 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; **25**: 3846–3852.
- 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; **100**: 207–212.
- Yerushalmi R, Woods R, Ravdin PM, Hayes MM, Gelmon KA. Ki67 in breast cancer: prognostic and predictive potential. *Lancet Oncol.* 2010; **11**: 174–183.
- 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; **103**: 1656–1664.
- 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; **153**: 477–491.
- Robertson JFR, Dowsett M, Bliss JMet al. Peri-operative aromatase inhibitor treatment in determining or predicting long term outcome in early breast cancer – the POETIC trial. San Antonio Breast Cancer Symposium, 6 December 2017; abstract GS1-03.
- Criscitello C, Disalvatore D, De Laurentis 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; **23**: 69–75.
- 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; **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; **29**: 1302–1308.
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; **19**: 296–300.
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; **8**: e019635.
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; **12**: 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; **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; **99**: 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; **35**: 1061–1069.
21. Polley MY, Leung SC, McShane LM *et al.* An international Ki67 reproducibility study. *J. Natl Cancer Inst.* 2013; **105**: 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; **8**: 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; **13**: e0188983.
24. Reinert T, Goncalves R, Ellis MJ. Current status of neoadjuvant endocrine therapy in early stage breast cancer. *Curr. Treat. Options Oncol.* 2018; **19**: 23.
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; **57**: 83–89.
26. Focke CM, Burger H, van Diest PJ *et al.* Interlaboratory variability of Ki67 staining in breast cancer. *Eur. J. Cancer* 2017; **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; **198**: 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; **54**: 1040–1048.
29. Polley MY, Leung SC, Gao D *et al.* An international study to increase concordance in Ki67 scoring. *Mod. Pathol.* 2015; **28**: 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; **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; **63**: 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; **48**: 787–794.
33. R Core Team. *R: a language and environment for statistical computing.* Vienna, Austria: R Foundation for Statistical Computing, 2018.
34. Vyberg M, Møller J, Røge R. Nordic immunohistochemical Quality Control – Ki67 assessment. 2018. Available at: <http://www.nordiqc.org/epitope.php?xml:id=1>. Last accessed 2019-06-13.
35. Acs B, Pelekanou V, Bai Y *et al.* Ki67 reproducibility using digital image analysis: an inter-platform and inter-operator study. *Lab. Invest.* 2019; **99**: 107–117.
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; **72**: 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; **169**: 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.* 2019; **32**: 59–69.

Supporting Information

Additional Supporting Information may be found in the online version of this article:

Figure S1. Specimen selection and study design flowchart.

Figure S2. Scoring field allocation algorithm.

Figure S3. Typewriter nuclei counting pattern.

Figure S4. Variability in Ki67 scores.

Figure S5. Percent agreement on categories.

Table S1. Cohort characteristics (from pathology reports and case notes) of the 30 cases.

Table S2. Summary statistics for weighted global scores (0–100%),¹ ordered according to observer median.

Table S3. Summary statistics for unweighted global scores (0–100%),¹ ordered according to observer median.

Table S4. Summary statistics for hot-spot scores (0–100%),¹ ordered according to observer median.

Table S5. Variance components estimates¹ and corresponding credible intervals.

Data S1. Scoring protocol.