

## Extent, impact, and mitigation of batch effects in tumor biomarker studies using tissue microarrays

Konrad H. Stopsack,<sup>1,2</sup> Svitlana Tyekucheva,<sup>3,4</sup> Molin Wang,<sup>2,3,5</sup> Travis A. Gerke,<sup>6</sup> J. Bailey Vasselkiv,<sup>2</sup> Kathryn L. Penney,<sup>2,5</sup> Philip W. Kantoff,<sup>1</sup> Stephen P. Finn,<sup>7,8</sup> Michelangelo Fiorentino,<sup>2,9</sup> Massimo Loda,<sup>10</sup> Tamara L. Lotan,<sup>11</sup> Giovanni Parmigiani,<sup>3,5#</sup> Lorelei A. Mucci<sup>2#</sup>

# joint senior authors

- <sup>1</sup> Department of Medicine, Memorial Sloan Kettering Cancer Center, New York, NY
- <sup>2</sup> Department of Epidemiology, Harvard T.H. Chan School of Public Health, Boston, MA
- <sup>3</sup> Department of Biostatistics, Harvard T.H. Chan School of Public Health, Boston, MA,
- <sup>4</sup> Department of Data Science, Dana-Farber Cancer Institute, Boston, MA, USA
- <sup>5</sup> Channing Division of Network Medicine, Department of Medicine, Brigham and Women's Hospital and Harvard Medical School, Boston, MA
- <sup>6</sup> Department of Cancer Epidemiology, Moffitt Cancer Center, Tampa, FL
- <sup>7</sup> Department of Pathology, St. James's Hospital, Dublin, Ireland
- <sup>8</sup> Trinity College, Dublin, Ireland
- <sup>9</sup> Pathology Unit, Addarii Institute, S. Orsola-Malpighi Hospital, Bologna, Italy
- <sup>10</sup> Department of Pathology, Weill Cornell Medical College, New York, NY
- <sup>11</sup> Department of Pathology, Johns Hopkins Medical Institutions, Baltimore, MD

*Correspondence:* Konrad H. Stopsack, Department of Medicine, Memorial Sloan Kettering Cancer Center, New York, NY, [stopsack@mskcc.org](mailto:stopsack@mskcc.org)

*Keywords:* tissue microarray; batch effects; measurement error; batchtma R package

*Running head:* Batch effects between tissue microarrays

1 **Abstract**

2 Tissue microarrays (TMAs) have been used in thousands of cancer biomarker studies. To what extent batch effects,  
3 measurement error in biomarker levels between slides, affects TMA-based studies has not been assessed  
4 systematically. We evaluated 20 protein biomarkers on 14 TMAs with prospectively collected tumor tissue from  
5 1,448 primary prostate cancers. In half of the biomarkers, more than 10% of biomarker variance was attributable to  
6 between-TMA differences (range, 1–48%). We implemented different methods to mitigate batch effects (R package  
7 *batchtma*), tested in plasmode simulation. Biomarker levels were more similar between mitigation approaches  
8 compared to uncorrected values. For some biomarkers, associations with clinical features changed substantially after  
9 addressing batch effects. Batch effects and resulting bias are not an error of an individual study but an inherent feature  
10 of TMA-based protein biomarker studies. They always need to be considered during study design and addressed  
11 analytically in studies using more than one TMA.

## 12 Introduction

13 Tissue microarrays (TMAs) were first developed in the 1990s as an efficient way to examine tissue-based  
14 biomarkers (1). Since then, TMAs have been used in thousands of studies to evaluate histologic and molecular  
15 biomarkers, mostly in cancer tissue. Even when biomarker assays are well standardized and run conditions are  
16 diligently kept fixed, some TMA slides (batches) may have measurements systematically too low or too high, and  
17 some batches may have wider spread around the true values of the biomarker than others. In general, such batch  
18 effects can have a profound impact on the validity of biomarker studies, such those using RNA microarrays (2, 3).  
19 Contrary to popular belief, whether such measurement error induces upward or downward bias in results is not  
20 guaranteed to follow simple heuristics (4).

21 Whether and to what extent TMAs are affected by batch effects has not been empirically assessed. TMAs  
22 pose unique challenges. For example, when tumor tissue is collected prospectively for inclusion on TMAs, tumor  
23 characteristics may differ between batches due to nonrandom assignment of cases, as well as temporal trends in tumor  
24 risk factors, screening, and diagnosis. Differences in tissue processing or storage across tissue specimens may have  
25 differential impact on biomarkers. Including calibration samples for quality control is also more challenging for TMAs  
26 than, for example, assaying of blood samples, because repeat sections from a tumor may differ due to intratumoral  
27 heterogeneity rather than only batch effects.

28 In this study, we assess batch effects in a large set of centrally constructed TMAs from prostate cancer tissue  
29 from 1,448 men in two nationwide cohort studies. We quantify the extent to which protein biomarker variation could  
30 be explained by batch effects. We probe different methods for mitigating batch effects while maintaining true,  
31 “biological,” between-TMA variation, including in a plasmode simulation. Finally, we demonstrate the impact of  
32 handling batch effects on commonly performed biomarker analyses.

## 33 Results

34  
35 **Extent and type of batch effects.** To evaluate the presence of batch effects in studies using TMAs, we studied tumor  
36 tissue from 1,448 men with primary prostate cancer on 14 TMAs, each including multiple tumor cores from 47 to 158  
37 patients per TMA (Figure 1). Multiple cores from the same tumor (usually 3) were always located on the same TMA.

38 TMAs were used to quantify 20 protein biomarkers (Figure 2). Biomarker values showed noticeable between-  
39 TMA variation. We estimated that across the 20 biomarkers, between-TMA variation explained between 1% and 48%  
40 of overall variation in biomarker levels (intraclass correlation coefficient, ICC), with half of the biomarkers having  
41 ICCs greater than 10% (Figure 2).

42 In an example biomarker, estrogen receptor alpha in nuclei of stromal cells (Figure 3), the means of the most  
43 extreme TMAs differed by 2.2 standard deviations in intensity of expression and variances differed up to 9.3-fold.  
44 Other biomarkers showed similar between-TMA variation by magnitude and by which TMAs had the most extreme  
45 values (Figure 4A). Likewise, we observed that not only means, but also variances of biomarker levels differed  
46 between TMAs, although patterns of heteroskedasticity appeared weaker than for means (Suppl. Figure 1). In contrast,  
47 we found little evidence for more complex patterns of batch effects, such that tumors with specific grade, stage, or  
48 year of diagnosis would have been particularly affected by between-TMA differences (Suppl. Table 1). Nevertheless,  
49 observations from the same TMAs tended to be clustered together when projected onto the first two principal  
50 components, capturing 27% of variance in all biomarkers (Figure 4B).

51 The method of scoring, including human (eye) scoring and computer-assisted quantification, differed between  
52 biomarkers, as did the main quantitative score, typically a measure of staining intensity, a proportion of cells above an  
53 intensity threshold, or a combination of both (Figure 2). Notably, between-TMA differences were present with any of  
54 these approaches. For example, batch effects were not only present when considering intensities of biomarker staining,  
55 as for the estrogen receptor alpha and beta example. Even when setting cut-offs for staining visible by eye and  
56 quantifying the number of stain-positive cells, 8% (95% CI, 2 to 15) of variance in estrogen receptor alpha positivity  
57 and 27% (95% CI, 11 to 42) of estrogen receptor beta positivity were attributable to between-TMA variation (Suppl.  
58 Figure 2).

59 In summary, we observed a large and concerning degree of between-TMA variation for several biomarkers  
60 that were quantified using different approaches, suggesting that addressing batch effects could significantly impact  
61 scientific inference.

62  
63 **Source of batch effects.** The noticeable proportion of variance attributable to TMAs could have two possibly co-  
64 existing explanations. First, that between-TMA differences in biomarkers reflect different patient and tumor  
65 characteristics that need to be retained. Second, that between-TMA differences are artifacts due to systematic  
66 measurement error that need to be removed (batch effects).

67 In support of the first hypothesis, there were noticeable differences in patient and tumor characteristics  
68 between TMAs that are likely associated with biomarker levels (Figure 1). Along with a 14-year range between the  
69 per-TMA medians of cancer diagnosis year, there were differences in the proportion of tumors with a Gleason score of  
70 8 or higher (between 11% and 33%) and rates of lethal disease (between 2 and 16 events per 1000 person-years of  
71 follow-up).

72 In support of the second hypothesis, we observed that certain TMAs had consistently higher or lower  
73 biomarker values for the majority of tested biomarkers (Figure 4A). For example, the same batches that showed  
74 higher-than-average biomarker values for stathmin also had higher-than-average values for PTEN. This example is  
75 noteworthy because both markers were assayed together on the same section of each TMA using multiplex  
76 immunofluorescence, and stathmin would be expected to be expressed in more aggressive tumors with activation of  
77 the PI3K signaling pathway while PTEN expression would be expected to be low in the same tumors (5).

78 Further supporting the second hypothesis, we did not observe any meaningful reduction in ICCs when we  
79 considered tumors that had the same clinical features in terms of Gleason score and stage (Suppl. Figure 3). Moreover,  
80 the association between Gleason score and biomarker levels (Figure 2D) was considerably lower than between TMAs  
81 and biomarker levels, as underscored by less pronounced visual separation of principal components by Gleason score  
82 (Figure 4C) than by TMA (Figure 4B). Gleason score differences explained no more than 13% of variance in  
83 biomarker levels (for prostate-specific membrane antigen, PSMA; 95% CI for ICC, 0.02 to 0.29), and 13 of the 20  
84 biomarkers had ICCs by Gleason score of 1% or less (Suppl. Figure 4).

85 To directly disentangle both hypotheses, we further examined data on 10 tumors with a total of 53 tumor  
86 cores for which some cores were included on different TMAs (Figure 4D). These were not included in the previous  
87 analyses and had estrogen receptor scoring data. This design allowed us to estimate biomarker differences directly  
88 attributable to between-TMA variability within the same tumors while controlling for the between-core variability  
89 expected due to intratumoral heterogeneity. Of the total variance in estrogen receptor alpha levels, 30% (95% CI, 0 to  
90 67) was explained by between-TMA variation; for estrogen receptor beta, 24% (95% CI, 0 to 60) was explained by  
91 between-TMA variation. For comparison, between-tumor variation explained 37% (95% CI, 4 to 68) of the variance  
92 of estrogen receptor alpha levels and 26% (95% CI, 0 to 57) of the variance of estrogen receptor beta levels.

93 Collectively, while these observations highlighted moderate differences in clinical and pathological  
94 characteristics between TMAs, they suggested that between-TMA differences were largely due to batch effects.

95  
96 **Mitigation of batch effects.** We implemented six different approaches for batch effects mitigation and compared  
97 these to the uncorrected biomarker levels (Figure 3, Suppl. Figure 5). Two mitigation approaches, batch means  
98 (approach 2) and quantile normalization (approach 6), assumed no true difference between TMAs is arising from  
99 patient and tumor characteristics, while all other approaches attempted to retain such differences between TMAs.  
100 Overall, correlations between values adjusted by different approaches were higher (mean Pearson  $r$ , 0.97 to 1.00) than  
101 between uncorrected values and corrected values ( $r$ , 0.90 to 0.95), regardless of mitigation approach (Figure 4E).

102 Approaches 2–7 reduced visible separation by batch on plots of the first two principal components (Suppl.  
103 Figure 6). Variance attributable to between-TMA differences decreased to ICCs of <1% for all markers (Suppl.  
104 Table 2). An exception was the quantile regression-based approach 5; the ICCs after this approach remained up to  
105 10%. This method does not explicitly address differences in means between batches but allows associations between  
106 clinical and pathological factors and biomarker levels to differ at high and low quantiles (Suppl. Figure 7).

107 The differences between uncorrected values and batch effect-corrected values were remarkably similar  
108 between the mean-based approaches using approaches 2 (simple means), 3 (standardized batch means), and 4 (inverse  
109 probability-weighted batch means; Suppl. Figure 8). Consequently, batch effect-corrected values by approaches 2–4  
110 were highly correlated (Figure 4E). All mean-only batch effect mitigations also gave the same results when fitting  
111 outcome models stratified by batch (Suppl. Figure 9). However, batch-specific results differed for approaches that  
112 targeted between-batch differences in the variance of biomarkers.

113  
114 **Validating batch effect mitigation in plasmode simulation.** To compare the performance of the different batch  
115 mitigation approaches in a time-to-event analysis, we applied plasmode simulation (6) to fix the expected strength of  
116 the biomarker exposure–outcome relationship *a priori* before artificially introducing batch effects. The correlation  
117 structure between biomarker and confounders and between confounders and batches from the actual data (Suppl.  
118 Figure 10A, C) was preserved in the plasmode-simulated data. Likewise, across a range of hazard ratios for the  
119 biomarker–outcome association, confounder–outcome associations remained unchanged (Suppl. Figure 10B, D).

120 We first evaluated a setting in which we did not introduce batch effects (Figure 5A). Here, the observed  
121 hazard ratios without batch effect mitigation equaled the expected. When performing (unnecessary) batch effect  
122 mitigation, observed hazard ratios were still comparable with the expected hazard ratios (Figure 5D; see Suppl.  
123 Table 3 for confidence intervals).

124 We then introduced batch effects by adding batch-specific mean differences to the observed biomarker levels,  
125 yet without introducing differences in variance by batch (Figure 5B). Without batch effect mitigation, for a true hazard  
126 ratio of 3.0, the observed hazard ratio, averaged over simulations, was 2.17 (95% CI, 1.86 to 2.53), an underestimate  
127 by 28% (Figure 5E; Suppl. Table 3). In contrast, all mitigation approaches produced CIs that covered the expected  
128 hazard ratio (*e.g.*, approach 6 quantile normalization: hazard ratio, 3.03; 95% CI, 2.48 to 3.69).

129 When we introduced batch-specific differences in both means and in variances (Figure 5C), the observed  
130 hazard ratio without batch effect mitigation decreased to 1.90 (95% CI, 1.66 to 2.16) compared to the expected hazard  
131 ratio of 3.0 (Figure 5F; Suppl. Table 3). Batch effect mitigation methods that only focus on means (approaches 2–4)  
132 reduced but did not fully eliminate bias, with hazard ratios ranging between 2.67 and 2.70. Methods that address  
133 differences in both mean and variance resulted in less bias, with an observed hazard ratio of 3.11 (95% CI, 2.54 to  
134 3.81) for approach 6 (quantile normalization).

135 We also included two stratification-based approaches. Fitting survival models separately by batch, followed  
136 by inverse-variance pooling (approach 8) resulted in approximately unbiased estimates but was less efficient than  
137 other approaches, comes with a risk of sparse-data bias, and resulted in considerably wider confidence intervals in our  
138 simulation. Including batch as a stratification variable in a single Cox model (approach 9) was unbiased and efficient.  
139 A drawback of both stratification-based approaches is that they do not explicitly estimate batch effect-adjusted  
140 biomarker values that could be visualized directly.

141 Scenarios evaluated thus far were based on the actual, modest imbalance of confounders between batches and  
142 at most weak associations between the biomarker and confounders, resulting in weak confounding overall. We  
143 additionally introduced both modest and strong associations between biomarker and confounders and created more  
144 severe imbalance between batches (Suppl. Figure 11). In all scenarios, the ranking of mitigation methods was  
145 preserved (Suppl. Figure 12, Suppl. Tables 3 and 4), with the least bias obtained through quantile normalization  
146 (approach 6). Bias occurred when using uncorrected biomarker levels in the presence of any batch effects, except if  
147 there was no association between biomarker and outcome (*i.e.*, a hazard ratio of 1), and with mean-only approaches 2–  
148 4 if variance was also affected by batch effects. In no situation, except possibly with the quantile regression-based  
149 approach 5, were estimates after batch effect mitigation farther from the expected values than results based on  
150 uncorrected biomarker levels.

151  
152 **Impact of batch effects.** To illustrate how batch effect mitigations alter the results of commonly conducted tumor  
153 biomarker analyses, we estimated how uncorrected and corrected biomarker levels were associated with Gleason  
154 score and with rates of lethal disease. For markers with little between-TMA variability (low ICCs) such as beta-  
155 catenin, there were no noticeable differences in associations between using unadjusted and adjusted biomarker levels

156 irrespective of adjustment model, as expected from plasmode simulation. However, for markers with higher between-  
157 TMA variability (higher ICC) and stronger associations with the outcome, adjustment approaches led to noticeable  
158 differences (Figure 6). For example, uncorrected stathmin expression levels were not associated Gleason score  
159 (difference, 0.00 standard deviations per 1 grade-group increase; 95% CI, -0.05 to 0.05), while the difference in levels  
160 corrected according to approach 6 was 0.04 (95% CI, 0.00 to 0.07), suggesting a potentially qualitatively different  
161 interpretation (Figure 6A; Suppl. Table 5). In models for lethal disease (Figure 6B), the otherwise unadjusted hazard  
162 ratio for the highest quartile of the vitamin D receptor, compared to the lowest quartile, was 0.44 (95% CI, 0.23 to  
163 0.86); after mitigation using approach 6, the hazard ratio was 0.19 (95% CI, 0.09 to 0.40), suggesting that unadjusted  
164 biomarker levels could underestimate the prognostic association by 2.3-fold (Suppl. Table 6 and 7).

## 165 Discussion

166  
167 The key strength of using TMAs is their utility in parallelizing the assessment of biomarkers on a large number of  
168 tissue specimens (1). Similar to other high-throughput platforms, batch effects have to be considered in every TMA  
169 biomarker study. As we demonstrated, for some of the biomarkers, batch effects can be of substantial magnitude. We  
170 show that batch effect mitigation is possible and can enhance study findings.

171 In our study of prostate tumor specimens, between-TMA differences explained 10% or more of the variance  
172 in biomarker levels for half of the included biomarkers, considerably more than one of the strongest pathological  
173 features in prostate cancer, Gleason grade. All analytical mitigation approaches to reduce batch effects, whether they  
174 attempted to retain real differences between tumors from different TMAs or not, led to corrected biomarker levels that  
175 were more similar to each other than they were, in general, to the uncorrected biomarker levels. In drawing from a  
176 large set of protein tumors biomarkers in prostate cancer, we show how appropriately mitigating batch effects  
177 strengthens results and their validity for biomarkers affected by batch effects.

178 Ideally, batch effects between TMAs are minimized when designing a study. Standardizing how tumor  
179 samples are obtained, stored, processed, and assayed is critical, as are stratified or random allocation of tumor samples  
180 to different TMAs (2) when possible. However, the batch effects that we observed occurred despite all feasible  
181 standardization efforts. Moreover, samples will be collected sequentially, and TMAs may be constructed sequentially  
182 in large-scale prospective studies over time. There were modest differences in the clinical and pathological  
183 characteristics between our TMAs, an issue that may be inevitable in larger-scale biobank studies. Allocation schemes  
184 of tumors to TMAs that appear ideal retrospectively, for example by matching “cases” of lethal tumors with  
185 “controls” of non-lethal tumors, may not be feasible prospectively. Likewise, in few of the thousands of studies using  
186 TMAs will it be possible to reallocate tumors to different TMAs and repeat all pathology work merely to reduce  
187 implications of batch effects.

188 An additional challenge in the design phase is that tissue samples are inherently heterogeneous and cannot  
189 simply be diluted, like blood samples. “Quality control” tumor samples that could serve as a quantitative calibration  
190 series suitable for all future biomarkers do not exist. One potential strategy is to include cell lines that have been  
191 formalin-fixed and paraffin-embedded on each TMA. While cell lines address issues of heterogeneity, the cell lines  
192 are often genomically unique and as such may not be relevant for all biomarkers. Another potential approach is to  
193 include samples from the same tumor case across TMAs, which would allow for direct estimation of batch effects. For  
194 these reasons, a principled approach that anticipates batch effects and addresses them analytically is critical.

195 Beyond efforts to prevent batch effects during the study design phase, we suggest the following best practices  
196 when undertaking TMA-based tissue biomarker studies (Figure 7). First, the extent of potential batch effects should be  
197 explored and reported in any study of cancer tissue using TMAs. Inspecting TMA slides and plots (Figure 3) (7) is  
198 important. Between-TMA variation should be quantified, for example by calculating ICCs, *i.e.*, to contrast variation of  
199 biomarker levels between TMAs compared to that between or within tumors (8). In our study, for half of the  
200 biomarkers, ICCs for between-TMA variation were low, at less than 10%, although the proportion of tolerable batch  
201 variation should be chosen based on the context. Whether TMAs differ in terms of average biomarker levels, low  
202 levels (possibly reflective of background), or variability between tumors will also inform what impact of between-  
203 TMA differences to expect.



204 Second, the source of between-TMA differences should be elucidated. Ideally, including multiple cores from  
205 the same tumors in more than one TMA will help estimating, again using ICCs, how biomarker levels vary between  
206 TMAs, between tumors, and within tumors. Alternatively, ICCs between TMAs can be estimated by restricting to or  
207 adjusting for tumor features associated with differences in the biomarker, if known. In our study, both approaches  
208 indicated that the largest share of between-TMA differences was likely due to batch effects rather than due to true  
209 differences between tumors on different TMAs. In multidisciplinary team discussions (9), it may be possible to  
210 directly pinpoint the source of batch effects and eliminate its cause. For example, if immunohistochemical staining  
211 was performed separately for each TMA, then immunohistochemistry and quantification should be repeated using new  
212 sections from all TMAs at once. Imaging of pathology slides can also be a source of batch effects (10). In other cases,  
213 particularly if such obvious reasons for batch effects were avoided through standardized processing, as in our  
214 examples, it may remain elusive whether batch effects were induced through subtle differences in how tumors were  
215 cored and embedded during TMA construction, how long they had been stored, how they were sectioned, how well  
216 the staining process was standardized, or how successfully background signal was eliminated during software-based  
217 quantification. Yet even biomarkers scored by manual quantification were not free from batch effects.

218 Third, if a biomarker is affected by batch effects and no “physical” remediation is possible, then analytical  
219 approaches should be used to reduce bias in results (2, 3). We demonstrate that in all plausible or exaggerated real-  
220 world scenarios, estimates after applying batch effect mitigations were consistently closer to the true underlying  
221 values than they were without. If batches do not only differ in terms of mean values, but also in terms of their  
222 variances, then methods that focus solely on means may be insufficient. A simple quantile-normalization-based  
223 approach was successful in reducing bias in real-world scenarios and could be preferred for its simplicity. It is  
224 important to note that any method tested in this study is preferable over not addressing batch effects, and thus the  
225 choice between methods should be secondary to the choice to address batch effects altogether. Only results for  
226 biomarkers that are affected by batch effects and that are associated with the outcome of interest will show large  
227 changes in estimates, as the vitamin D receptor in our example. In contrast, for the majority of our example  
228 biomarkers, results did not change appreciably because batch effects were low, associations with the outcome were  
229 close to null, or both (Figure 6).

230 We recommend that researchers openly address batch effects in their TMA-based studies: they are not an  
231 error of an individual study, but an inherent feature of TMA-based studies. Batch effects have long been recognized in  
232 studies of the transcriptome using microarrays and next-generation sequencing, where batch effect mitigations are a  
233 component of standard workflows (3, 11). Our data strongly suggest that protein biomarker studies using multiple  
234 TMAs are at risk of batch effects just like any other biomarker study. The extent of batch effects is difficult to predict,  
235 and empirical evaluation is necessary each time. Future studies should quantify between-TMA differences and, if they  
236 deem batch effect mitigations to be unnecessary, provide evidence for absence of batch effects, rather than merely  
237 assuming their absence. The methods that we provide facilitate appropriate migration of batch effects between TMAs  
238 and help strengthen scientific inference. It may be prudent to extend this approach to in-situ tissue biomarkers other  
239 than proteins, such as RNA in-situ hybridization, even if our study only demonstrated batch effects for proteins.  
240 Having mitigated batch effects will allow researchers to focus on increasing study validity by addressing other sources  
241 of measurement error (4), selection bias (for example, from tumor biospecimen availability) (12), and confounding.

## 242 **Methods**

243  
244 **TMAs and biomarkers.** Tumor tissue in this study was from men who were diagnosed with primary prostate cancer  
245 during prospective follow-up of two nationwide cohort studies. The Health Professionals Follow-up Study is an  
246 ongoing cohort study that enrolled 51,529 male health professionals across the United States in 1986. The Physicians’  
247 Health Study 1 and 2 were randomized-controlled trials of aspirin and dietary supplements, starting in 1982 with  
248 22,071 male physicians. Participants were diagnosed with and treated for prostate cancer at local health care providers  
249 across the United States. The study team collected formalin-fixed paraffin-embedded tissue specimens from radical  
250 prostatectomy and transurethral resection of the prostate (TURP), and study genitourinary pathologists performed  
251 central re-review, including standardized Gleason grading of full hematoxylin–eosin-stained slides from the tumor

252 blocks (13). The study protocol was approved by the institutional review boards of the Brigham and Women's  
253 Hospital and Harvard T.H. Chan School of Public Health, and those of participating registries as required.

254 TMAs were constructed using 0.6-mm tissue cores of the primary nodule or the nodule with the highest  
255 Gleason score (14), including three or more cores of tumor tissue per participant (tumor). For a subset of tumors,  
256 additional cores of tumor-adjacent, histologically normal-appearing prostate tissue were included. TMAs were  
257 constructed at the same laboratory across a 10-year period, as tissue from cohort participants became available,  
258 without matching on patient or tumor characteristics and without randomization. Cores from the same participant  
259 were generally included on the same TMA, with exceptions noted below, and summarized as the mean. We include  
260 information from 14 prostate tumor tissue microarrays.

261 Immunostaining was generally performed separately for individual biomarkers yet always for all TMAs at the  
262 same time. Detailed immunohistochemistry staining and quantification procedures for each marker have been  
263 published (5, 15-25) or are in preparation for estrogen receptor alpha (antibody SP1; Thermo Scientific, Waltham,  
264 MA) and an antibody (PPG5/10; Bio-Rad Laboratories, Hercules, CA) widely used to measure estrogen receptor beta.  
265 If batch effect mitigation approaches had been applied in the original studies, the uncorrected levels were retrieved.  
266 Right-skewed biomarker scores (Ki-67, pS6, TUNEL) were  $\log_e$  transformed. All biomarkers were scaled to mean 0  
267 and standard deviation 1 solely to facilitate comparisons of batch effects across markers; batch effect mitigation does  
268 not necessitate scaling and preserves absolute biomarker values.

269 **Extent and type of batch effects.** To visualize the extent of biomarker variation between TMAs, we plotted  
270 uncorrected biomarker values by tumor, biomarker, and TMA. We summarized biomarker variation using the first two  
271 principal components. We calculated between-TMA mean differences and ratios of variances versus the first TMA.  
272 We tested if tumors with different clinical/pathological characteristics had higher biomarker levels in TMAs with  
273 higher means (*i.e.*, multiplicative effect modification). For each biomarker and each clinical/pathological feature  
274 (ordinal Gleason score, ordinal stage, or calendar year of diagnosis), let  $Z_{ij}$  be the within-TMA  $z$ -score (mean 0,  
275 standard deviation 1) for tumor  $i$  from TMA  $j$ ;  $A_i$ , the clinical/pathological feature of tumor  $i$ ;  $B_j$ , the TMA-specific  
276 biomarker mean,  $r_j$ , the TMA-specific random effect, and  $e_{ij}$ , residual error. In the regression model  
277  $Z_{ij} = \beta_0 + \beta_1 A_i + \beta_2 B_j + \beta_3 A_i B_j + r_j + e_{ij}$ , we evaluated the  $\beta_3$  term to assess for multiplicative effect measure  
278 modification.

280 We calculated the proportion of variation in biomarker levels attributable to TMA using intra-class  
281 correlations (ICCs, also "repeatability" (8)) based on one-way random effects linear mixed models with an  
282 independent variance-covariance structure (8, 26) for  $Y_{ij}$ , the biomarker level per tumor  $i$  and TMA  $j$ ; where  $\beta_0$  is the  
283 biomarker mean;  $r_j$ , the random effect for TMA  $j$ ; and  $e_{ij}$ , the residual error:  $Y_{ij} = \beta_0 + r_j + e_{ij}$ . The ICC was defined  
284 as the proportion of between-TMA variance in the total variance:  $ICC = \frac{\text{var}(r)}{\text{var}(r) + \text{var}(e)}$ . 95% CIs for ICCs were  
285 obtained using parametric bootstrapping using 500 repeats (27).

286 **Source of batch effects.** To directly distinguish between-TMA variation caused by batch effects from variation  
287 caused by differences in patient and tumor characteristics, we compared ICCs per biomarker overall to ICCs per  
288 biomarker when restricting analyses to a subset of tumors with the same clinical features. We also leveraged a small  
289 subset of tumors that had cores included on more than one TMA. Here, we used two-way random effects linear mixed  
290 models with independent variance-covariance structure to separate between-TMA variation from between-core  
291 variation (*i.e.*, intratumoral heterogeneity) and residual modeling error:  $Y_{ijk} = \beta_0 + r_j + s_i + e_{ijk}$ . Compared to the  
292 model described earlier, this model additionally includes tumor-specific random effects  $s_i$ , and thus  
293

$$294 \quad ICC = \frac{\text{var}(r)}{\text{var}(r) + \text{var}(s) + \text{var}(e)}.$$



296 **Mitigation of batch effects.** In addition to (1) using uncorrected values, we implemented eight different approaches  
 297 to handle between-TMA batch effects:

298 (2) *Simple means.* This approach assumes that all TMAs, if not affected by batch effects, would have the  
 299 same mean biomarker value. Differences in mean biomarker values per batch are corrected by estimating batch-  
 300 specific mean effects (differences from the overall mean level) using a linear regression model with uncorrected  
 301 biomarker values as the outcome and batch indicators as predictors. Corrected biomarker values are then obtained by  
 302 subtracting batch-specific effects from the uncorrected biomarker values. Mean differences can either indicate the  
 303 difference of each batch mean to the overall mean, as implemented here, or be defined by comparison to a reference  
 304 batch.

305 (3) *Standardized means.* This approach estimates marginal means per batch using model-based  
 306 standardization (in the epidemiologic sense). It assumes that batches with similar characteristics have the same means  
 307 if not affected by batch effects. A linear regression model for a specific biomarker is fit, adjusting for tumor variables  
 308 that differ in distribution between TMAs, similar to an approach described in the epidemiology literature by  
 309 Rosner (28). Let  $Y_{ij}$  indicate the biomarker value for tumor  $i$  on TMA  $j$ ;  $B_j$ , TMA  $j$ ;  $C_1$  to  $C_m$ , the  $m$  covariates to be  
 310 retained; and  $e_{ij}$ , the residuals. Then  $Y_{ij} = \beta_0 + \beta_j B_j + \gamma_1 C_1 + \dots + \gamma_n C_n + e_{ij}$ . Batch effect-corrected biomarker  
 311 values can be obtained by subtracting batch-specific effects  $\beta_j$  predicted from the model above from uncorrected  
 312 biomarker values.

313 We included the following clinical and pathologic variables as plausible sources of real between-TMA  
 314 differences that should be retained in this approach, as well approaches 4–7: calendar year of diagnosis (linear),  
 315 Gleason score (categorical: 5–6; 3+4; 4+3; 8; 9–10), and pathologic tumor stage (categorical: pT1/T2, pT3/T3a,  
 316 pT3b/T4/N1, missing/tissue from transurethral resection of the prostate).

317 (4) *Inverse-probability weighted batch means.* Like the preceding approach, this approach assumes that  
 318 batches with similar characteristics have the same means if not affected by batch effects. While the preceding  
 319 approach assumes a constant association between covariates and biomarker levels across batches, this approach allows  
 320 for associations to differ between batches. We first used inverse probability weighting for marginal standardization of  
 321 the distribution of clinical and pathological features per batch to the distribution in the entire study population.  
 322 Stabilized weights (29), truncated at the 2.5<sup>th</sup> and 97.5<sup>th</sup> percentile, were obtained through multinomial regression  
 323 models, modeling the probability of assignment to a specific batch based on same clinical and pathological variables  
 324 as in (3). In the weighted pseudopopulation, we then used a marginal linear model to estimate batch-specific mean  
 325 differences, which were further used as in approaches 2 and 3.

326 (5) *Quantile regression.* This approach assumes that batches with similar characteristics have the same values  
 327 for a selected set of batch-specific quantiles, in this application the upper and lower quartile. The lower quartile may  
 328 be particularly affected by background noise, while the upper quartile may more likely reflect differences in batches  
 329 due to covariates. A corollary of separately modeling the two differently is that clinical and pathological variables are  
 330 allowed to have different effects on these quartiles (30). These assumptions contrast with approaches 2–4 that focus on  
 331 mean levels only. We used quantile regression with the Frisch-Newton approach separately for the first and third  
 332 quartile of biomarker values with batch indicators to predict adjusted batch-specific quantile values with the same  
 333 confounders as above. We then used the batch-specific 25<sup>th</sup> percentiles ( $y^{\tau=0.25}$ ) as the offset and the interquartile  
 334 range between the 25<sup>th</sup> and 75<sup>th</sup> percentiles ( $y^{\tau=0.75}$ ) as the scaling factor when batch-correcting biomarker levels. Let  
 335  $y_{ij}^*$  indicate the batch effect-corrected biomarker level for tumor  $i$  on TMA  $j$ ;  $y_{ij}$ , the uncorrected biomarker level for  
 336 tumor  $i$  on TMA  $j$ ;  $\hat{y}_j^{\tau=x}$ ,  $x^{\text{th}}$  quantile of  $y$  for batch  $j$  (predicted value for  $y_j$  from unadjusted quantile regression);  
 337  $\hat{y}_j^{\tau=x,*}$  is  $\hat{y}_j^{\tau=x}$  with adjustment for confounders (predicted value for  $y_j$  from adjusted quantile regression); and  $\bar{y}^{\tau=x}$ ,  
 338 the  $x^{\text{th}}$  quantile of  $y$  overall. Then the corrected biomarker level is

$$339 y_{ij}^* = \frac{(y_{ij} - \hat{y}_j^{\tau=0.25}) (\bar{y}^{\tau=0.75} - \bar{y}^{\tau=0.25})}{(\hat{y}_j^{\tau=0.75,*} - \hat{y}_j^{\tau=0.25,*})} + \bar{y}^{\tau=0.25} - \hat{y}_j^{\tau=0.25,*} + \hat{y}_j^{\tau=0.25}$$

340 (6) *Quantile normalization.* This approach assumes that samples on all batches, if not affected by batch  
 341 effects, would not only have the same mean and variance but also the same distribution of individual biomarker

342 values. Uncorrected biomarker values are ranked within each batch and then ranks are replaced by the mean of values  
343 with the same rank across batches. We implemented quantile normalization using *limma* (31).

344 A conceptually related approach, for example employed in molecular epidemiology (2, 9), would be to use  
345 within-batch ranks as the batch-corrected biomarker, often grouped into data-driven categories such as batch-specific  
346 quartiles. We did not further consider these derivatives because they do not retain absolute biomarker levels and can  
347 distort rank distances.

348 (7) *ComBat*. For comparison, we additionally included the ComBat algorithm, which like approach 4  
349 attempts to retain differences in batch means due to covariate differences; it is frequently applied together with  
350 approach 6. ComBat and its derivatives (11, 32, 33) were initially designed for microarray studies of gene expression,  
351 which include considerably more than one biomarker per sample. This property would typically limit their use for a  
352 protein biomarker quantified on a TMA unless a large number of biomarkers is available, as in our study. Mitigation  
353 depends on values of other biomarkers on the same batches. Even if multiple protein biomarkers were available, the  
354 non-randomly selected set of concomitantly available biomarkers may influence how batch effects are corrected.  
355 ComBat scales means and (optionally) variances while (optionally) retaining adjustment variables. ComBat is  
356 implemented using an empirical Bayes approach to achieve more favorable properties for small batches. The  
357 underlying model is similar to the regression above and has been emulated by a two-way analysis of variance (34). In  
358 using ComBat, we scaled both means and variances, adjusting for the same clinical and pathological variables as  
359 before. Because ComBat cannot handle biomarkers if they are missing on entire batches, we ran ComBat separately  
360 for groups of biomarkers measured on 8, 9, 10, or 14 TMAs.

361 (8) *Stratification with inverse-variance pooling*. An alternative approach to treating batch effects is to  
362 estimate outcome regression models separately by batch. This approach can be applied for a variety of regression  
363 models but does not result in corrected values. We pooled estimates with inverse variance-weighting to obtain  
364 summary estimates.

365 (9) *Stratification in Cox proportional hazards regression*. In a special case of stratification for time-to-event  
366 outcomes, Cox proportional hazards models allow for nonparametric batch effect mitigation by including batch as a  
367 stratification factor in the model specification. Comparisons are performed within batches. Unlike approach 8, only  
368 batch-specific baseline hazard functions but no batch-specific effects are estimated.

369 For approaches 1–7, we calculated Pearson correlation coefficients between uncorrected and corrected  
370 biomarker levels. Additionally, we repeated ICC and principal components analyses with corrected levels, and we  
371 estimated associations between Gleason score and biomarker levels after batch effect mitigation, stratifying by batch  
372 using approach 8.

373 Approaches 2–6, which result in batch effect-adjusted biomarker levels, are implemented in the R package  
374 *batchtma*, available at <https://stopsack.github.io/batchtma>.

375  
376 **Plasmode simulation.** We evaluated the impact of batch effect mitigation approaches on known, investigator-  
377 determined biomarker–outcome associations using plasmode simulation, an approach used, for example, for  
378 evaluating confounding control for binary exposures in pharmacoepidemiology (6). We used observed data from all  
379 tumors included on the 14 TMAs to determine covariates (Gleason grade, pathological stage) and outcome (lethal  
380 disease), preserving the observed correlation structure (*e.g.*, joint distribution of clinical characteristics across TMAs).  
381 The only simulated elements were the biomarker levels and the strengths of biomarker–outcome associations (hazard  
382 ratios ranging from  $\frac{1}{3}$  to 3) that we fixed by simulating event times with flexible parametric survival models (35).  
383 Models used a baseline hazard function consisting of cubic splines with three knots. Group differences were based on  
384 proportional hazards for the observed confounder–outcome coefficients in the real data and the fixed biomarker  
385 (exposure)–outcome hazard ratios.

386 First, we used plasmode simulation to generate the fixed associations of the biomarker levels with the  
387 outcome, which are unknowable outside simulation studies, generating 200 plasmode datasets for each association.  
388 Second, we introduced batch effects. Batch effects were either only for the mean or for both mean and variance, using  
389 the actual standardized between-TMA differences in means and variances for the estrogen receptor-alpha protein, a

390 biomarker with high ICCs. We also added batch effects for mean and variance with effect modification, making mean  
391 and variance changes due to batch effects strongly correlated with Gleason scores. Third, we calculated batch effect-  
392 adjusted biomarker levels using approaches 2–6. Lastly, we compared the expected hazard ratios for the biomarker-  
393 outcome association, fixed during simulations, with the estimated hazard ratios (with normality-based 95% CIs)  
394 before and after batch effect mitigation approaches 2–6 and using the two stratification-based approaches 8 and 9.

395 In sensitivity analyses, we simulated “moderate” associations between the biomarker and confounders  
396 (0.2 standard deviations difference in biomarker levels per Gleason grade group, 0.1 per stage category), “strong”  
397 associations (differences of 0.4 and 0.2 standard deviations, respectively; stronger than observed for any biomarker in  
398 our study), as well as “strong” associations and additional imbalance in Gleason grade and stage between TMAs (by  
399 excluding tumors with low grades from TMAs with higher-than-average means and excluding tumors with high stage  
400 from TMAs with low-than-average means), all before the four steps described above.

401  
402 **Impact of batch effects.** To quantify the impact of different approaches to batch-effect handling on scientific  
403 inference, we focused on two commonly employed types of analyses in biomarker research in prostate cancer: first, a  
404 cross-sectional analysis of Gleason score and biomarker levels, using linear regression models; second, a time-to-  
405 event analysis of biomarker levels and rates of lethal disease, using Cox proportional hazards regression. For  
406 graphing, exposures were modeled in five categories (Gleason scores) or using restricted cubic splines with three  
407 knots (all biomarkers in models for lethal disease). For numeric comparisons, Gleason scores were modeled as ordinal  
408 variables and biomarkers as linear variables to obtain one single estimate per model. We also categorized biomarkers  
409 into four quartiles and compared hazard ratios for lethal disease of the extreme quartiles. Models were designed only  
410 for investigating issues of batch effects and not for subject matter inference on specific biomarkers.

411  
412 **Data availability.** The batchtma R package is available at <https://stopsack.github.io/batchtma>. Code used to produce  
413 results this manuscript is at [https://github.com/stopsack/batchtma\\_manuscript](https://github.com/stopsack/batchtma_manuscript). Data are available for analysis on the  
414 Harvard FAS computing cluster through a project proposal for the Health Professionals Follow-up Study  
415 (<https://sites.sph.harvard.edu/hpfs/for-collaborators>).

## 416 **Acknowledgments**

417  
418 We thank the participants and staff of the HPFS and the PHS for their valuable contributions. In particular, we would  
419 like to recognize the contributions of Liza Gazeeva, Siobhan Saint-Surin, Robert Sheahan, Betsy Frost-Hawes, and  
420 Eleni Konstantis. We would like to thank the following state cancer registries for their help: AL, AZ, AR, CA, CO,  
421 CT, DE, FL, GA, ID, IL, IN, IA, KY, LA, ME, MD, MA, MI, NE, NH, NJ, NY, NC, ND, OH, OK, OR, PA, RI, SC,  
422 TN, TX, VA, WA, and WY. The authors assume full responsibility for analyses and interpretation of these data. The  
423 HPFS is supported by the NIH (U01 CA167552). This research was funded in part by the Specialized Programs of  
424 Research Excellence program in Prostate Cancer 5P50 CA090381 and P50 CA211024, the NIH/NCI Cancer Center  
425 Support Grants P30 CA008748 and P30 CA006516, NIH/NCI grants 5R37 CA227190-02 (S. Tyekucheva,  
426 K.L. Penney, G. Parmigiani, and L.A. Mucci), R03 CA212799 (M.W.), R35 CA212799 (M.W.), and R01 CA131945  
427 (M. Loda). The Department of Defense supported K.H. Stopsack (W81XWH-18-1-0330). K.H. Stopsack,  
428 K.L. Penney, S.P. Finn, and L.A. Mucci are Prostate Cancer Foundation Young Investigators.

## 429 **Competing Interests**

430  
431 P.W. Kantoff reports the following disclosures for the last 24-month period: he has investment interest in Context  
432 Therapeutics LLC, DRGT, Placon, and Seer Biosciences; he is a company board member for Context Therapeutics  
433 LLC; he is a consultant/scientific advisory board member for Bavarian Nordic Immunotherapeutics, DRGT, GE

434 Healthcare, Janssen, OncoCellMDX, Progenity, Seer Biosciences, and Tarveda Therapeutics; and he serves on data  
435 safety monitoring boards for Genentech/Roche and Merck.

436

437 G. Parmigiani reports the following disclosures for the last 24-month period: he had investment interest in CRA  
438 Health; he is a co-founder and company board member of Phaeno Biotechnology; he is a consultant / scientific  
439 advisory board member for Konica-Minolta, Delfi Diagnostics and Foundation Medicine; he serves on a data safety  
440 monitoring board for Geisinger. None of these activities are related to the content of this article.

## References

- 442 1. J. Kononen *et al.*, Tissue microarrays for high-throughput molecular profiling of tumor specimens. *Nat. Med.* **4**, 844-847  
443 (1998).
- 444 2. S. S. Tworoger, S. E. Hankinson, Use of biomarkers in epidemiologic studies: minimizing the influence of measurement error  
445 in the study design and analysis. *Cancer Causes Control* **17**, 889-899 (2006).
- 446 3. J. T. Leek *et al.*, Tackling the widespread and critical impact of batch effects in high-throughput data. *Nat Rev Genet* **11**, 733-  
447 739 (2010).
- 448 4. M. van Smeden, T. L. Lash, R. H. H. Groenwold, Reflection on modern methods: five myths about measurement error in  
449 epidemiological research. *Int. J. Epidemiol.* **49**, 338-347 (2020).
- 450 5. K. H. Stopsack *et al.*, Multiplex Immunofluorescence in Formalin-Fixed Paraffin-Embedded Tumor Tissue to Identify Single-  
451 Cell-Level PI3K Pathway Activation. *Clin. Cancer Res.* **26**, 5903-5913 (2020).
- 452 6. J. M. Franklin, S. Schneeweiss, J. M. Polinski, J. A. Rassen, Plasmode simulation for the evaluation of pharmacoepidemiologic  
453 methods in complex healthcare databases. *Comput Stat Data Anal* **72**, 219-226 (2014).
- 454 7. S. Manimaran *et al.*, BatchQC: interactive software for evaluating sample and batch effects in genomic data. *Bioinformatics* **32**,  
455 3836-3838 (2016).
- 456 8. S. Nakagawa, H. Schielzeth, Repeatability for Gaussian and non-Gaussian data: a practical guide for biologists. *Biol. Rev.*  
457 *Camb. Philos. Soc.* **85**, 935-956 (2010).
- 458 9. M. T. Marrone *et al.*, Adding the Team into T1 Translational Research: A Case Study of Multidisciplinary Team Science in the  
459 Evaluation of Biomarkers of Prostate Cancer Risk and Prognosis. *Clin. Chem.* **65**, 189-198 (2019).
- 460 10. S. Kothari *et al.*, Removing batch effects from histopathological images for enhanced cancer diagnosis. *IEEE J Biomed Health*  
461 *Inform* **18**, 765-772 (2014).
- 462 11. J. Leek, W. E. Johnson, A. Jaffe, H. Parker, J. Storey (2011) The SVA package for removing batch effects and other unwanted  
463 variation in high-throughput experiments.
- 464 12. L. Liu *et al.*, Utility of inverse probability weighting in molecular pathological epidemiology. *Eur. J. Epidemiol.* **33**, 381-392  
465 (2018).
- 466 13. J. R. Stark *et al.*, Gleason score and lethal prostate cancer: does  $3 + 4 = 4 + 3$ ? *J. Clin. Oncol.* **27**, 3459-3464 (2009).
- 467 14. A. Pettersson *et al.*, The TMPRSS2:ERG rearrangement, ERG expression, and prostate cancer outcomes: a cohort study and  
468 meta-analysis. *Cancer Epidemiol. Biomarkers Prev.* **21**, 1497-1509 (2012).
- 469 15. J. R. Rider *et al.*, Tumor expression of adiponectin receptor 2 and lethal prostate cancer. *Carcinogenesis* **36**, 639-647 (2015).
- 470 16. R. Flavin *et al.*, SPINK1 protein expression and prostate cancer progression. *Clin. Cancer Res.* **20**, 4904-4911 (2014).
- 471 17. M. Fiorentino *et al.*, Overexpression of fatty acid synthase is associated with palmitoylation of Wnt1 and cytoplasmic  
472 stabilization of beta-catenin in prostate cancer. *Lab. Invest.* **88**, 1340-1348 (2008).
- 473 18. T. U. Ahearn *et al.*, Calcium sensing receptor tumor expression and lethal prostate cancer progression. *J. Clin. Endocrinol.*  
474 *Metab.* 10.1210/jc.2016-1082, jc20161082 (2016).
- 475 19. Z. Ding *et al.*, SMAD4-dependent barrier constrains prostate cancer growth and metastatic progression. *Nature* **470**, 269-273  
476 (2011).
- 477 20. P. L. Nguyen *et al.*, Fatty acid synthase polymorphisms, tumor expression, body mass index, prostate cancer risk, and survival.  
478 *J. Clin. Oncol.* **28**, 3958-3964 (2010).
- 479 21. T. U. Ahearn *et al.*, Expression of IGF/insulin receptor in prostate cancer tissue and progression to lethal disease.  
480 *Carcinogenesis* **39**, 1431-1437 (2018).
- 481 22. A. Pettersson *et al.*, MYC Overexpression at the Protein and mRNA Level and Cancer Outcomes among Men Treated with  
482 Radical Prostatectomy for Prostate Cancer. *Cancer Epidemiol. Biomarkers Prev.* **27**, 201-207 (2018).
- 483 23. J. L. Kasperzyk *et al.*, Prostate-specific membrane antigen protein expression in tumor tissue and risk of lethal prostate cancer.  
484 *Cancer Epidemiol. Biomarkers Prev.* **22**, 2354-2363 (2013).
- 485 24. P. K. Dhillon *et al.*, Aberrant Cytoplasmic Expression of p63 and Prostate Cancer Mortality. *Cancer Epidemiology Biomarkers*  
486 *& Prevention* **18**, 595-600 (2009).
- 487 25. W. K. Hendrickson *et al.*, Vitamin D receptor protein expression in tumor tissue and prostate cancer progression. *J. Clin. Oncol.*  
488 **29**, 2378-2385 (2011).
- 489 26. S. E. Hankinson *et al.*, Reproducibility of plasma hormone levels in postmenopausal women over a 2-3-year period. *Cancer*  
490 *Epidemiol. Biomarkers Prev.* **4**, 649-654 (1995).
- 491 27. M. A. Stoffel, S. Nakagawa, H. Schielzeth, S. Goslee, rptR: repeatability estimation and variance decomposition by generalized  
492 linear mixed-effects models. *Methods Ecol. Evol.* **8**, 1639-1644 (2017).
- 493 28. B. Rosner, N. Cook, R. Portman, S. Daniels, B. Falkner, Determination of blood pressure percentiles in normal-weight children:  
494 some methodological issues. *Am. J. Epidemiol.* **167**, 653-666 (2008).
- 495 29. S. R. Cole, M. A. Hernan, Constructing inverse probability weights for marginal structural models. *Am. J. Epidemiol.* **168**, 656-  
496 664 (2008).
- 497 30. D. Bann, E. Fitzsimons, W. Johnson, Determinants of the population health distribution: an illustration examining body mass  
498 index. *Int. J. Epidemiol.* **49**, 731-737 (2020).
- 499 31. M. E. Ritchie *et al.*, limma powers differential expression analyses for RNA-sequencing and microarray studies. *Nucleic Acids*  
500 *Res.* **43**, e47 (2015).
- 501 32. W. E. Johnson, C. Li, A. Rabinovic, Adjusting batch effects in microarray expression data using empirical Bayes methods.  
502 *Biostatistics* **8**, 118-127 (2007).
- 503 33. Y. Zhang, D. F. Jenkins, S. Manimaran, W. E. Johnson, Alternative empirical Bayes models for adjusting for batch effects in  
504 genomic studies. *BMC Bioinformatics* **19**, 262 (2018).



- 505 34. V. Nygaard, E. A. Rodland, E. Hovig, Methods that remove batch effects while retaining group differences may lead to  
506 exaggerated confidence in downstream analyses. *Biostatistics* **17**, 29-39 (2016).  
507 35. M. J. Crowther, P. C. Lambert, Simulating biologically plausible complex survival data. *Stat. Med.* **32**, 4118-4134 (2013).  
508

509 **Figures**

510

511

**Figure 1. Characteristics of men with prostate cancer with tissue included on the 14 tumor tissue microarrays.**

512

**A**, Calendar years of cancer diagnosis, with thick lines indicating median, boxes interquartile ranges, and whiskers 1.5

513

times the interquartile range. **B**, Counts of tumors by Gleason score. **C**, Counts of tumors by pathological TNM stage

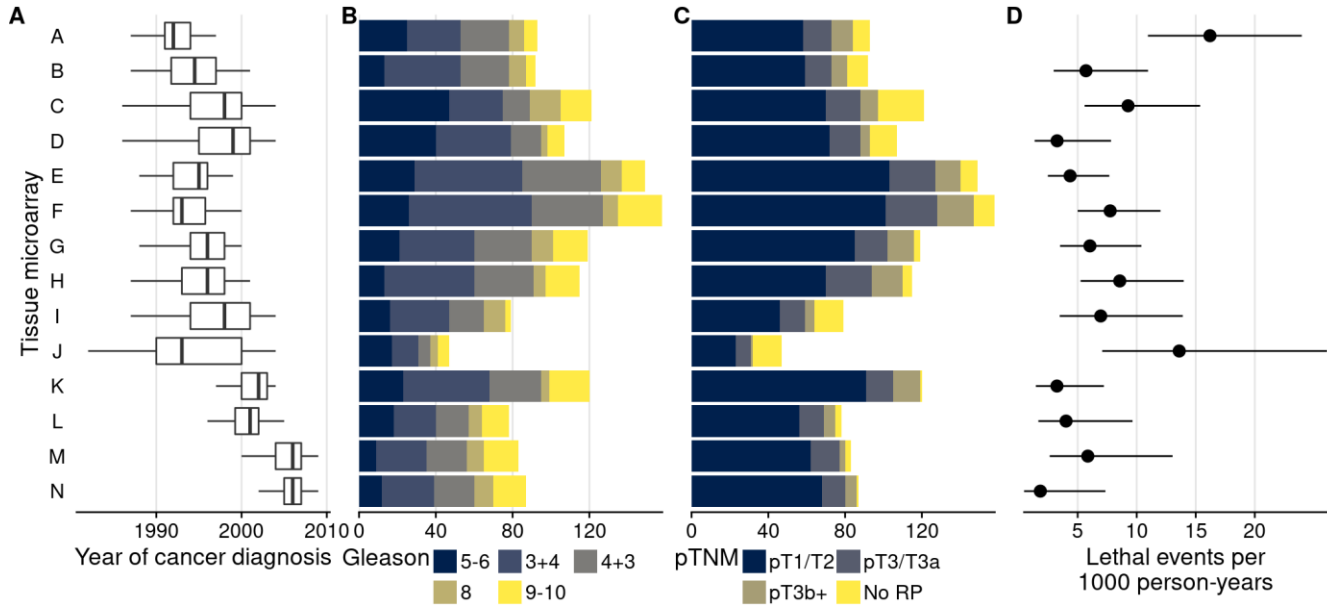
514

(RP: radical prostatectomy). **D**, Rates of lethal disease (metastases or prostate cancer-specific death over long-term

515

follow-up), with bars indicating 95% confidence intervals. As throughout, multiple cores are summarized per tumor.

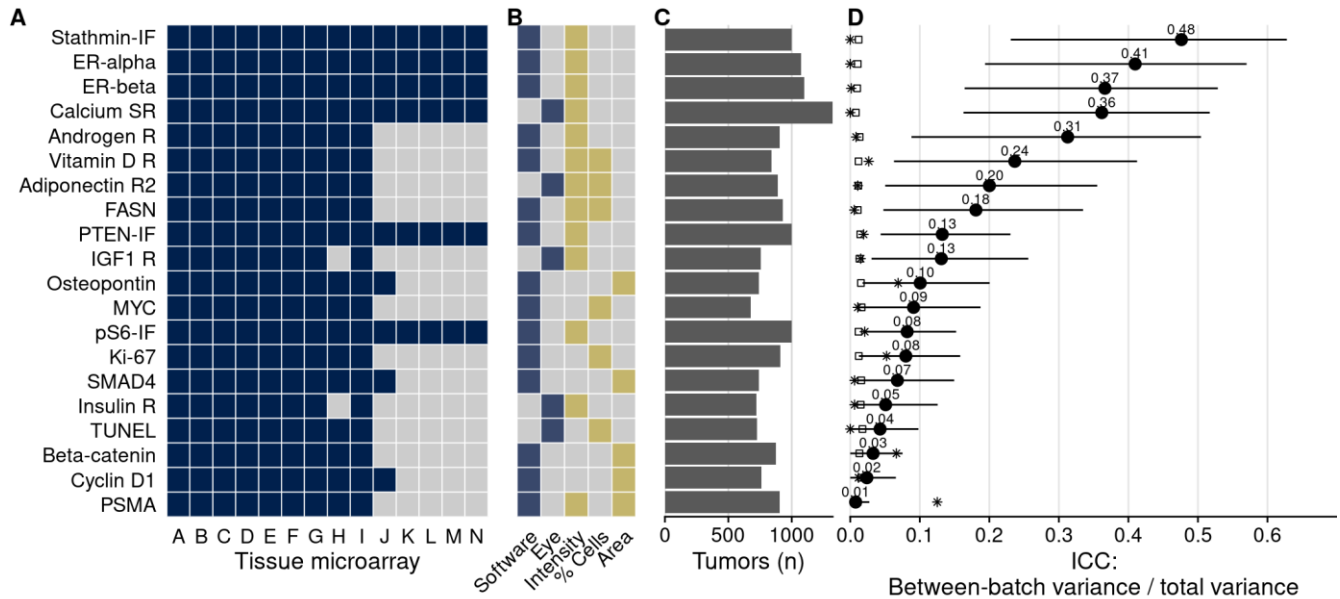
516



517

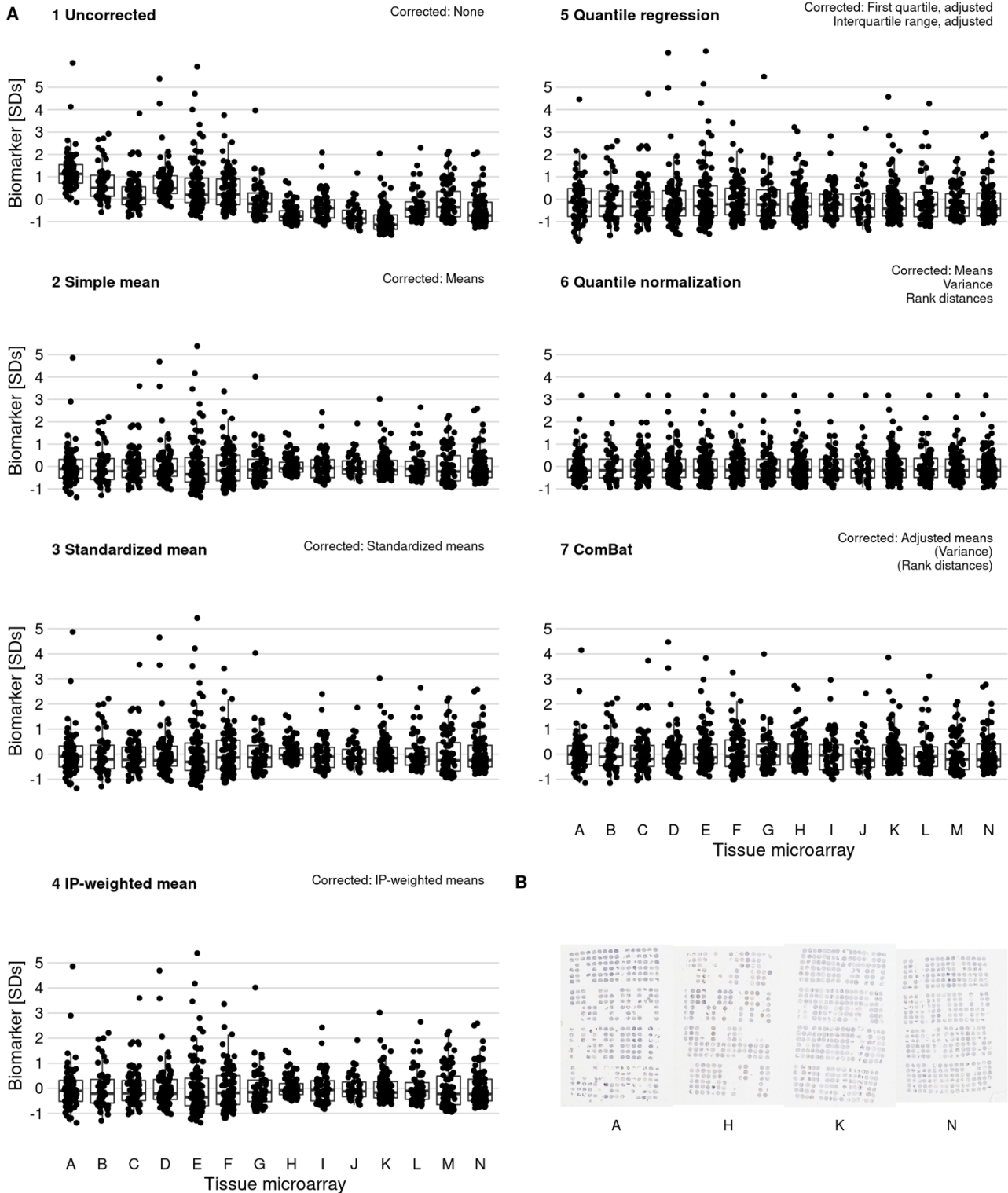
518 **Figure 2. Biomarkers stained, scoring methods, and intraclass correlation coefficients (ICCs).** **A**, Tissue  
 519 microarrays assessed for each marker (dark blue, yes). **B**, Approach to quantifying biomarkers: software-based  
 520 scoring vs. eye scoring (blue, yes); biomarker quality assessed: staining intensity, proportion of cells positive for the  
 521 biomarker, area of tissue positive for the biomarker (yellow, yes). **C**, Counts of tumors assessed for each biomarker.  
 522 **D**, Between-tissue microarray ICCs (*i.e.*, proportion of variance explained by between-tissue microarray differences)  
 523 for each biomarker, with 95% confidence intervals. Empty symbols indicate the 97.5<sup>th</sup> percentile of the null  
 524 distribution of the ICC obtained by permuting tumor assignments to TMAs; asterisks indicate between-Gleason grade  
 525 group ICCs. Biomarkers are arranged by descending between-tissue microarray ICC.

526



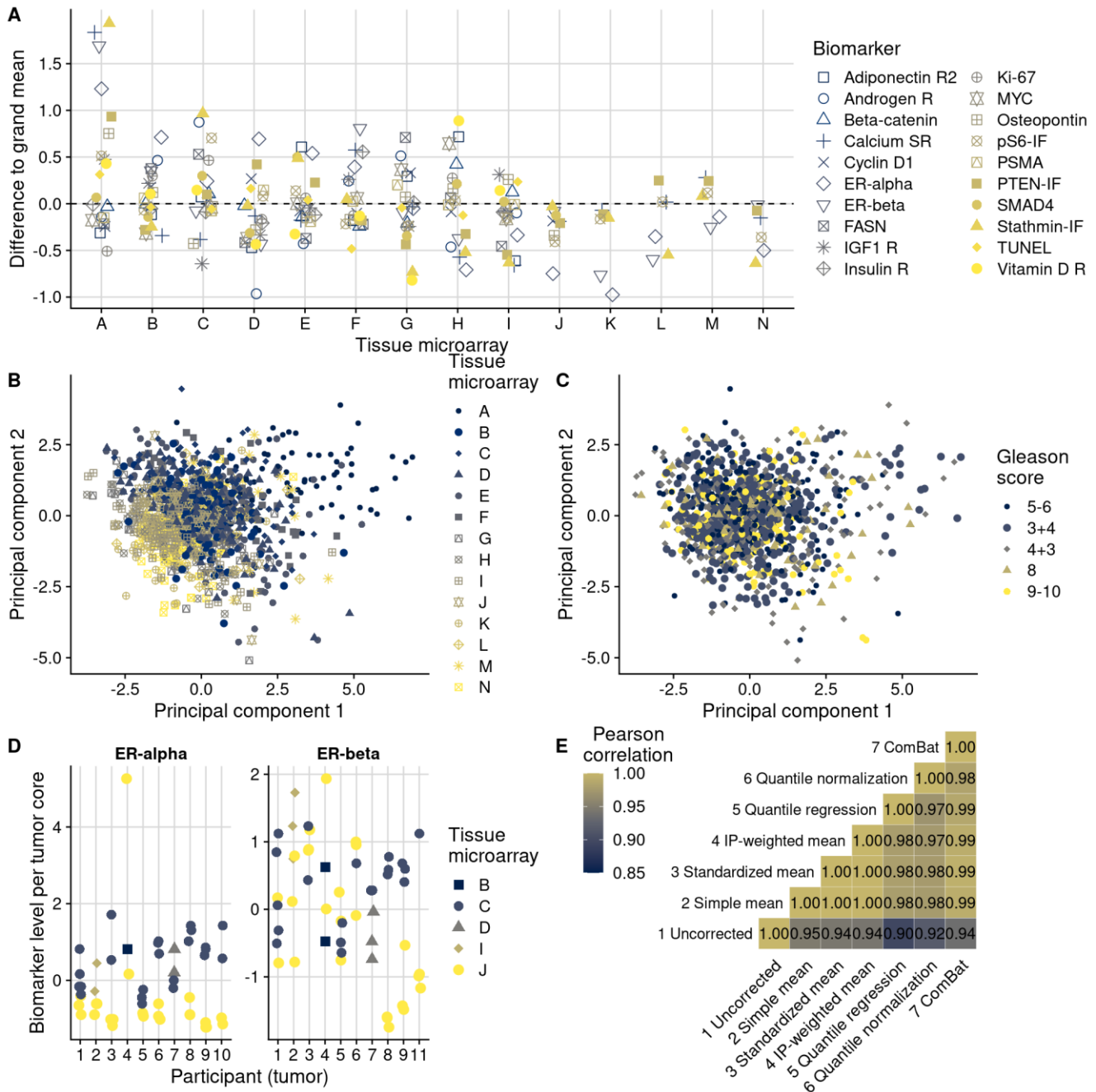
527

528 **Figure 3. Effect of batch effect mitigation on a biomarker with strong between-tissue microarray variation. A,**  
 529 The protein biomarker estrogen receptor- $\alpha$  was quantified as staining intensity in nuclei of epithelial cells,  
 530 averaged over all cores of each tumor. Each panel shows processed data for a specific approach to correcting batch  
 531 effects. Notes in the upper right corner indicate which properties of batch effects were potentially addressed. Each data  
 532 point represents one tumor.  $y$ -axes are standard deviations of the combined data for the specific method. Thick lines  
 533 indicate medians, boxes interquartile ranges, and whisker length is 1.5 times the interquartile range. **B,** Example  
 534 photographs of tissue microarrays; brown color indicates positive staining.



535

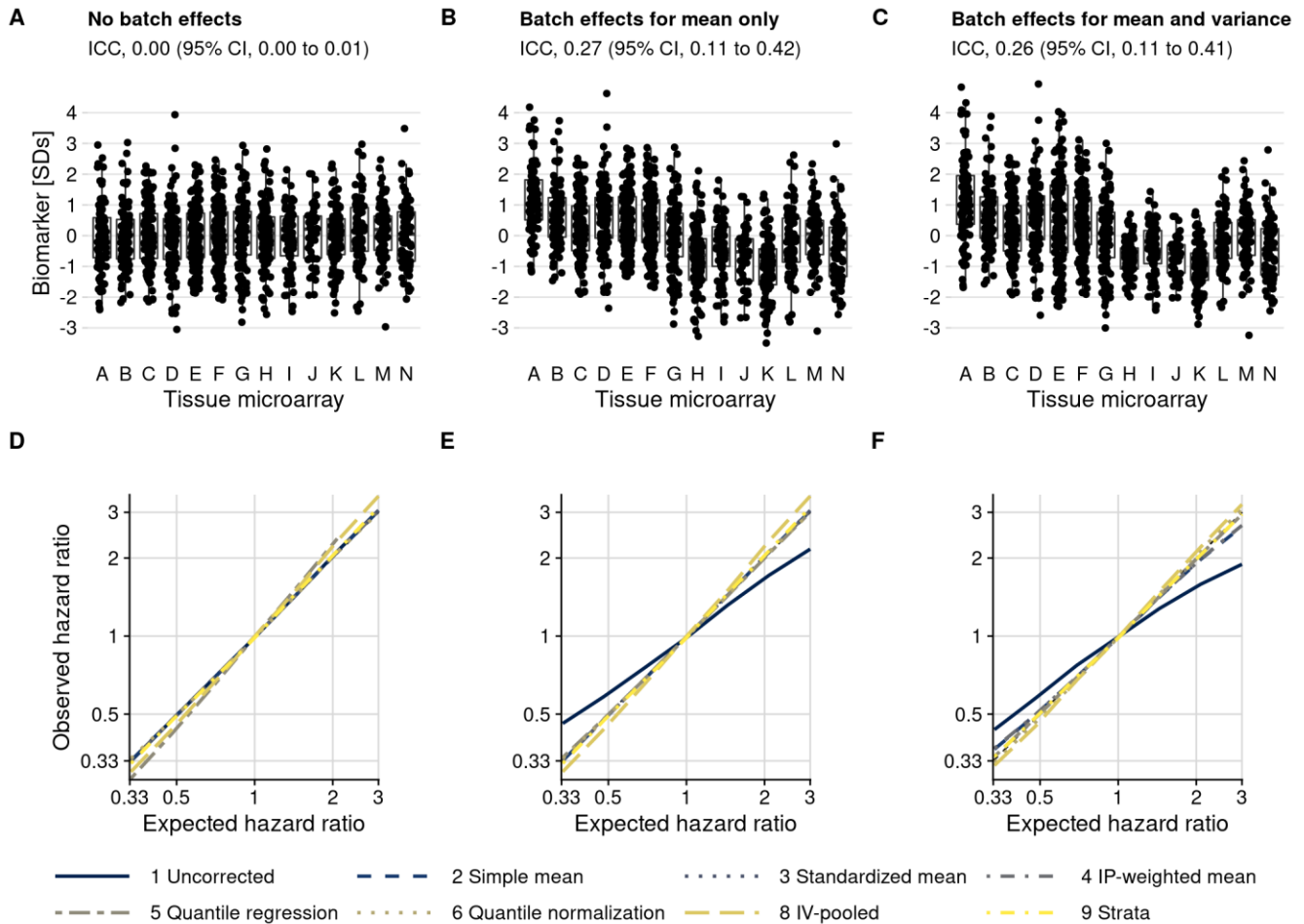
536 **Figure 4. Patterns, source, and remediation of batch effects.** **A**, Biomarker mean levels by tissue microarray, in  
 537 biomarker-specific standard deviations ( $\nu$ -axis). Each point is one tissue microarray. **B**, First two principal  
 538 components of biomarkers levels on all 14 tissue microarrays, with color/shape denoting tissue microarray. Each point  
 539 is one tumor. **C**, The same first two principal components, with color/shape denoting Gleason score. **D**, Per-core  
 540 biomarker levels for tumors with multiple cores included on two separate tissue microarrays, for estrogen receptor  
 541 (ER) alpha and beta, both in standard deviations. Each point is one tumor core. **E**, Pearson correlation coefficients  $r$   
 542 between uncorrected and corrected biomarker levels. Entries are averages across all markers.



543

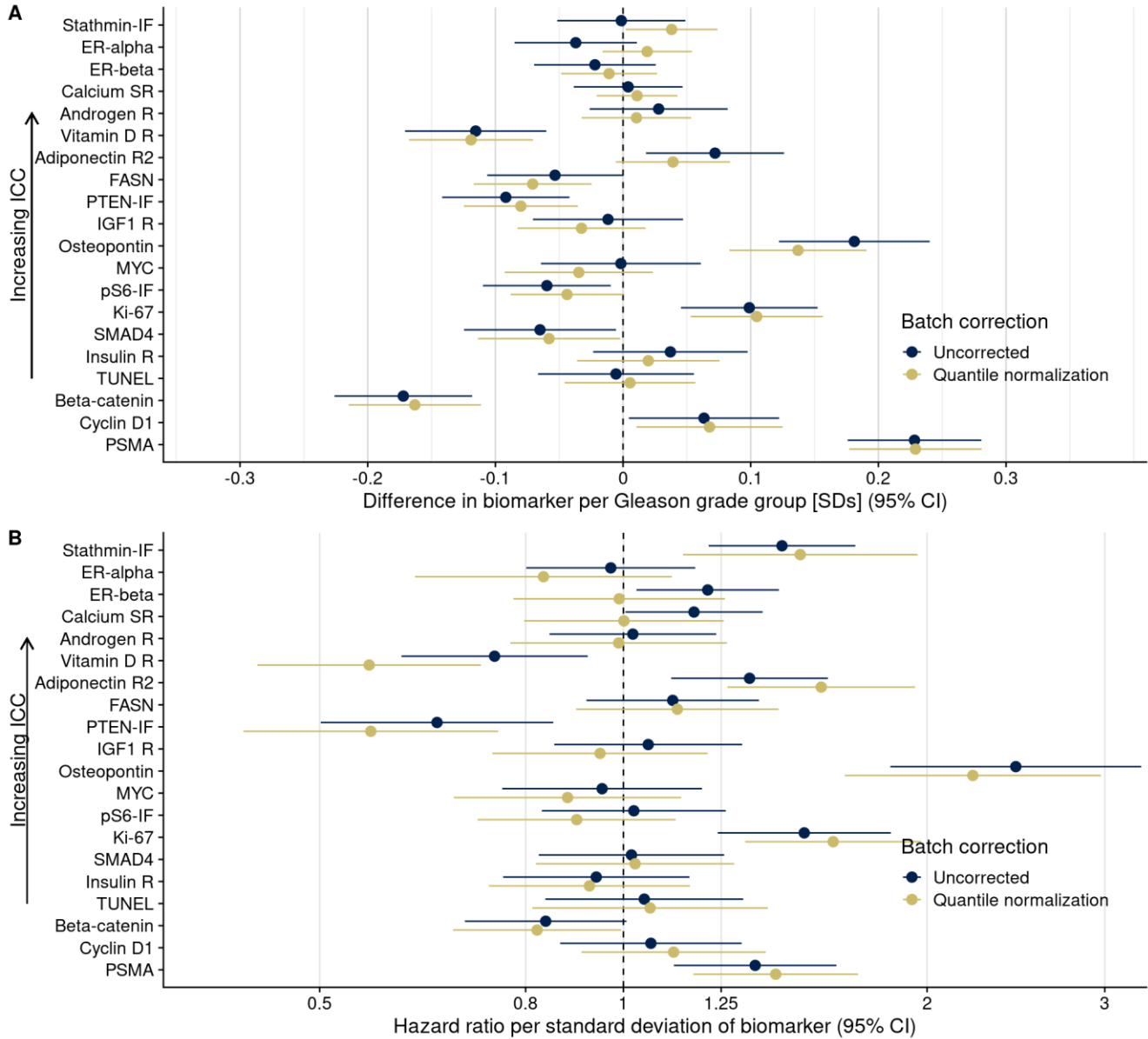


544 **Figure 5. Plasmode simulation results.** A–C, Biomarker levels by tissue microarray in three simulation scenarios;  
 545 **D–F**, true versus observed hazard ratios for the biomarker–outcome association after alternative approaches to batch  
 546 effect correction, with correction methods being numbered as in the Methods section. The solid blue line indicates no  
 547 correction for batch effects. **A** and **D**, no batch effects; **B** and **E**, means-only batch effects; **C** and **F**, means and  
 548 variance batch effects.



549

550 **Figure 6. Consequences of batch effect mitigation on scientific inference. A,** Gleason score and biomarker levels  
 551 (in standard deviations, per Gleason grade group). **B,** Biomarker levels and progression to lethal disease, with hazard  
 552 ratios per one standard deviation increase in biomarker levels from univariable Cox regression models. In both panels,  
 553 blue dots indicate estimates using uncorrected biomarker levels, yellow dots indicate batch effect-corrected levels,  
 554 applying approach (5), quantile regression. Lines are 95% confidence intervals. Biomarkers are ordered by decreasing  
 555 between-tissue microarray intraclass correlation coefficient (ICC).

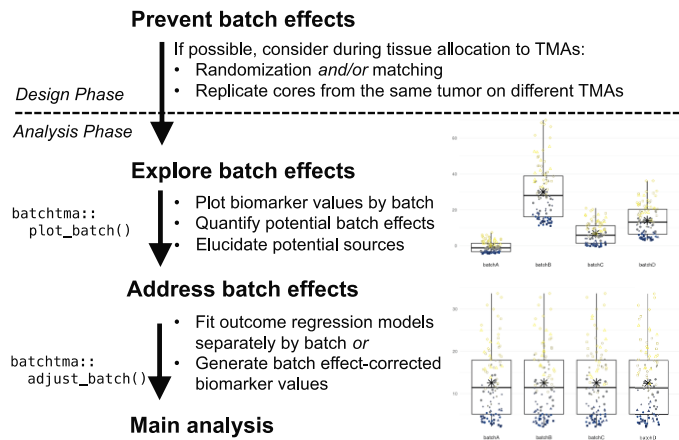


556

557  
558  
559  
560  
561  
562

### Figure 7. Recommended workflow for anticipating and handling batch effects between tissue microarrays.

Following prevention approaches at the design phase of a project, all tissue microarray-based studies should explore the potential for batch effects once a biomarker has been measured. Addressing batch effects should only be skipped there is no indication for their presence. Batch effect-corrected biomarker levels can easily be generated by the *batchtma* R package.



563

564 **Legends for Supplementary Tables and Figures**

565 See separate markdown document, also available at

566 [https://stopsack.github.io/batchtma\\_manuscript/batchtma\\_manuscript\\_210416.html](https://stopsack.github.io/batchtma_manuscript/batchtma_manuscript_210416.html)

567

568 **Supplementary Table 1.** Interaction terms to test for multiplicative effect modification, *i.e.* whether batch effects  
569 more strongly affect tumors with more extreme clinical/pathological characteristics. The table shows point estimates  
570 (differences in biomarker levels), 95% confidence interval bounds, p-values, and false-discovery rates (FDR, in  
571 ascending order) for interaction terms between the within-batch normalized biomarker level and the potential effect  
572 modifier in linear models that have absolute biomarker levels in standard deviation units per biomarker as the outcome  
573 and also include main effects for the biomarker and the effect modifier (terms not shown).

574 **Supplementary Table 2.** Intraclass correlation coefficient (ICC) for between-batch variance for uncorrected  
575 biomarker levels (“1 Raw”) and biomarker levels after applying different correction methods.

576 **Supplementary Table 3.** Results from plasmode simulation according to type of induced batch effect, using the data  
577 correlation structure “moderate confounding.” For three fixed (“true”) hazard ratios for the biomarker–outcome  
578 association ( $1/3$ , 1, and 3), the observed hazard ratios after batch correction (with 95% confidence intervals) are shown.

579 **Supplementary Table 4.** Results from plasmode simulation according to data correlation structure, using the batch  
580 effect “mean and variance.” For three fixed (“true”) hazard ratios for the biomarker–outcome association ( $1/3$ , 1, and  
581 3), the observed hazard ratios after batch correction (with 95% confidence intervals) are shown.

582 **Supplementary Table 5.** Gleason grade—biomarker associations according to batch effect correction method. Point  
583 estimates from unadjusted linear regression models for biomarker values with Gleason score categories per 1 “grade  
584 group” increase as the predictor are shown (with 95% confidence intervals). For  $\log_e$ -transformed markers like Ki-67,  
585 estimates are differences on the  $\log_e$  scale.

586 **Supplementary Table 6.** Biomarker levels and lethal disease according to batch effect correction method. Hazard  
587 ratios (with 95% confidence intervals) per 1 standard deviation increase in the biomarker (linear) from unadjusted Cox  
588 regression models are shown.

589 **Supplementary Table 7.** Biomarker levels and lethal disease according to batch effect correction method. Unlike in  
590 the preceding table, the hazard ratios (with 95% confidence intervals) are contrasts comparing extreme quartiles  
591 (fourth compared to first quartile) from unadjusted Cox regression models.

592

593 **Supplementary Figure 1.** Ratios of variance per tissue microarray to the mean variance for each marker.

594 **Supplementary Figure 2.** Tissue microarrays and differences in % positivity, at the example of estrogen receptor  
595 alpha and beta, and variance in biomarker levels explained by between-tissue microarray differences (ICC).

596 **Supplementary Figure 3.** Intraclass correlation coefficients (ICCs), quantifying the proportion of variance in  
597 biomarker levels attributable to between-tissue microarray differences, across all tumors and after restriction to those  
598 378 tumors across tissue microarrays that have the same clinical/pathological characteristics in terms of Gleason score  
599 3+4 and prostatectomy stage pT1/T2.

600 **Supplementary Figure 4.** Intraclass correlation coefficients (ICCs), quantifying the proportion of variance in  
601 biomarker levels attributable to between-Gleason grade differences, by increasing ICC.

602 **Supplementary Figure 5.** Uncorrected compared with batch effect-corrected biomarker levels, for estrogen receptor  
603 alpha. Symbols and color indicate the tissue microarray.

604 **Supplementary Figure 6.** Principal components 1 and 2 after batch effect correction using (5) quantile regression for  
605 biomarkers available on all tissue microarrays. Symbol color and shape indicate the tissue microarray.

606 **Supplementary Figure 7.** Quantile-specific associations of confounders (clinical/pathological differences) with  
607 (uncorrected) biomarker levels of estrogen receptor alpha. Shown are regression coefficients for the 10<sup>th</sup>, 50<sup>th</sup>, and 90<sup>th</sup>  
608 percentiles as the outcomes of quantile regression models.

609 **Supplementary Figure 8.** Batch corrections per tissue microarray and method. The plot shows the difference between  
610 uncorrected and corrected values per batch, averaged across all biomarkers. IGF1-R was excluded because of missing  
611 values for some correction approaches. For batch correction approaches that only address the mean (*i.e.*, that subtract  
612 the same correction value from all biomarker values within each batch), only that difference is visible; for correction  
613 methods that address individual values within batches differently, batch-specific medians and interquartile ranges of  
614 differences between uncorrected and corrected values are visible.

615 **Supplementary Figure 9.** Biomarker differences, after batch effect correction methods, for a one-unit increment in  
616 Gleason score, stratified by tissue microarray. “Pooled” indicates estimates pooled over batches (TMAs) using  
617 inverse-variance weighting. “No stratification” indicates estimates without stratification. Note that for batch effect  
618 correction approaches that only address between-batch differences in means (approaches 2–4), estimates stratified by  
619 batch (and pooled estimates thereof) are the same.

620 **Supplementary Figure 10.** Data structures in the actual data and in 200 plasmode simulation datasets. **A**, Gleason  
621 scores and lethal prostate cancer (metastasis-free survival) in the actual data. **B**, Gleason scores and lethal prostate  
622 cancer in an example simulated dataset. Shaded areas indicate 95% confidence intervals. **C**, Pearson correlation  
623 coefficients between biomarker levels and confounders, and between confounders, across all simulated datasets.  
624 Correlation coefficients observed in the actual data are noted in the legend. **D**, Hazard ratios for the biomarker and the  
625 confounders in relation to lethal prostate cancer, pooling all simulated data sets. Confounder–outcome associations  
626 were simulated to correspond to their observed values in the actual data; exposure–outcome associations were  
627 simulated to a range of hazard ratios (*x* axis). Lines indicate medians across simulations with the same exposure–  
628 outcome hazard ratio, shaded areas range from the 2.5<sup>th</sup> to 97.5<sup>th</sup> percentile.

629 **Supplementary Figure 11.** The data correlation structure “confounding and imbalance.” Tumors with more extreme  
630 Gleason scores were set to be more extremely influenced by batch effects in terms of mean and variances.

631 **Supplementary Figure 12.** Plasmode simulation results for all scenarios. Observed hazard ratios after different  
632 approaches to batch effect correction (*y* axis) are compared to true (fixed) hazard ratios for the biomarker–outcome  
633 association (*x* axis; solid blue line: no correction for batch effects). Columns are different batch effects that were  
634 added; rows are different data correlation structures.