1 **The prognostic effects of somatic mutations in ER-positive breast cancer**
2

3 Authors:
4 Obi L Griffith, PhD[1,2,3,4,*], Nicholas C Spies, BSc[1,*], Meenakshi Anurag, PhD[5,6*], Malachi Griffith,
5 PhD[1,2,3,4], Jingqin Luo, PhD[3,6], Dongsheng Tu, PhD[8], Belinda Yeo, PhD[9], Jason Kunisaki, BSc[1],
6 Christopher A Miller, PhD[1,2], Kilannin Krysiak, PhD[1,2], Jasreet Hundal, MSc[1], Benjamin J
7 Ainscough, BSc[1], Zachary L Skidmore, MEng[1], Katie Campbell, BSc[1], Runjun Kumar, BSc[2],
8 Catrina Fronick, BSc[1], Lisa Cook, BSc[1], Jacqueline E Snider, BSc[2], Sherri Davies, PhD[2], Shyam
9 M Kavuri, PhD[5,6], Eric C Chang, PhD[5,6], Vincent Magrini, PhD[1,4,10], David E Larson, PhD[1],
10 Robert S Fulton, MSc[1,4], Shuzhen Liu, MSc[8], Samuel Leung, MSc[8], David Voduc, MD[8], Ron
11 Bose, MD, PhD[2], Mitch Dowsett PhD, FMedSci[9], Richard K Wilson, PhD[1,3,4], Torsten O Nielsen,
12 MD, PhD[8], Elaine R Mardis, PhD[1,3,4,10,†], Matthew J Ellis MB, BChir, PhD[5,6†]
13

14 Affiliations:
15   1. McDonnell Genome Institute, Washington University School of Medicine, St. Louis, MO
16   2. Department of Medicine, Division of Oncology, Washington University School of
17      Medicine, St. Louis, MO
18   3. Siteman Cancer Center, Washington University School of Medicine, St. Louis, MO
19   4. Department of Genetics, Washington University School of Medicine, St. Louis, MO
20   5. Lester and Sue Smith Breast Center and Dan L. Duncan Cancer Center, Baylor College
21      of Medicine, Houston, TX
22   6. Department of Medicine, Baylor College of Medicine, Houston, TX
23   7. Division of Biostatistics, Washington University School of Medicine, St. Louis MO
24   8. Genetic Pathology Evaluation Centre, University of British Columbia, Vancouver,
25      Canada
26   9. Institute of Cancer Research, London, UK
27   10. Current address: Nationwide Children's Hospital and Department of Pediatrics, The Ohio
28      State University College of Medicine, Columbus, OH
29

30 * These authors contributed equally.
31 † Corresponding authors. mjellis@bcm.edu, elaine.mardis@nationwidechildrens.org
32

33

34 Author emails:
35 obigriffith@wustl.edu, nspies@wustl.edu, anurag@bcm.edu, mgriffit@wustl.edu,
36 jingqinluo@wustl.edu, dtu@ctg.queensu.ca, belinda.yeo@onjcri.org.au, kunisakijh@wustl.edu,
37 c.a.miller@wustl.edu, kkrysiak@wustl.edu, jhundal@wustl.edu, b.ainscough@wustl.edu,
38 zskidmor@wustl.edu, katiecampbell@wustl.edu, runjunkumar@wustl.edu, cfronick@wustl.edu,
39 cooklisa@wustl.edu, jsnider@dom.wustl.edu, sdavies@dom.wustl.edu,
40 meghashyam.kavuri@bcm.edu, echang1@bcm.edu, vincent.magrini@nationwidechildrens.org,
41 delarson@wustl.edu, rfulton22@wustl.edu, shuzhensuzanne.liu@vch.ca,
42 samuel.leung@vch.ca, dvoduc@bccancer.bc.ca, rbose@dom.wustl.edu,
43 mitchell.dowsett@icr.ac.uk, richard.wilson@nationwidechildrens.org, torsten@mail.ubc.ca,
44 elaine.mardis@nationwidechildrens.org, matthew.ellis@bcm.edu

52    Abstract
53
54    More than 50 genes are recurrently affected by somatic mutation in estrogen receptor positive
55    (ER+) breast cancer but prognostic effects have not been definitively established. Primary tumor
56    DNA was therefore subjected to targeted sequencing from 625 postmenopausal (UBC-TAM
57    series) and 328 premenopausal (MA12 trial) hormone receptor-positive (HR+) patients.
58    Independent validation of prognostic interactions was achieved using independent data from the
59    METABRIC study. Associations between MAP3K1 and PIK3CA with luminal A status and TP53
60    mutations with Luminal B/non-luminal tumors were observed, validating the methodological
61    approach.  In UBC-TAM, *NF1* frame-shift nonsense *(FS/NS)* mutation was validated as a poor
62    outcome driver. For MA12, poor outcome associated with PIK3R1 mutation was similarly
63    validated. DDR1 mutations were strongly associated with poor prognosis in UBC-TAM despite
64    stringent false-discovery correction (q=0.0003). In conclusion, uncommon recurrent somatic
65    mutations should be further explored to create a more complete explanation of the highly
66    variable outcomes that typify ER+ breast cancer.
67
68    Introduction
69
70    While recent genomic studies have provided a comprehensive catalog of genes that accumulate
71    somatic point mutations and small insertions/deletions (indels) in estrogen receptor-positive
72    (ER+) breast cancer, there remains considerable uncertainty as to how these newly discovered
73    mutations relate to disease outcomes[1-3]. Most genomic discovery cohorts were neither uniformly
74    treated nor followed long enough. For ER+ disease in particular, prognostic studies require
75    prolonged observation since late relapses can occur[4]. Uniform treatment was a feature of a
76    whole genome sequencing study of samples accrued from a neoadjuvant aromatase inhibitor
77    (AI) clinical trial for ER+ clinical stage 2 or 3 disease, although only short-term anti-proliferative
78    response to AI was reported. This investigation identified that mutations in *MAP3K1*, a tumor
79    suppressor gene involved in stress kinase activation, were associated with indolent biological
80    features and low proliferation rates[5]. The resulting hypothesis was that *MAP3K1* mutation would
81    be associated with favorable outcomes.  In contrast, *TP53* mutations associated with poor
82    prognosis features and high proliferation rates.
83
84    To more comprehensively address the relationships between somatic mutations and outcomes
85    in ER+ breast cancer, we developed an approach to detect somatic mutations in DNA isolated
86    from formalin fixed tumor blocks that were over 20 years old. After curating existing mutational
87    data from breast cancer genomics discovery studies (Supplementary Data 1), 83 genes were
88    chosen for analysis (Supplementary Table 1). We applied DNA hybrid capture, sequencing and
89    somatic analysis to three ER+ breast cancer discovery cohorts with contrasting clinical
90    characteristics: An older cohort treated with adjuvant tamoxifen and no chemotherapy, a
91    premenopausal cohort uniformly treated with chemotherapy and randomized to tamoxifen
92    versus observation; and a third mixed cohort that was used to expand the mutational landscape
93    analysis (Supplementary Table 2). An analytical pipeline was developed to identify somatic
94    variants while compensating for the lack of matched normal DNA, which is generally unavailable
95    in the setting of older formalin-fixed tumor material.  Somatic mutations were analyzed for
96    association with standard clinical variables, wherein mutated *TP53* and *MAP3K1* served as *a*
97    *priori* hypotheses for poor and good outcome, respectively. Additional objectives were to identify
98    new mutational hotspots and to determine mutation frequencies for therapeutic targets.
99    Validation was possible by comparing our results to those in cBioPortal where the mutational
100   analysis in the METABRIC cohort overlapped with the 83 genes investigated in the study
101   described here.
102

103    Results
104
105    Sequencing and final study cohorts
106
107    University of British Columbia Tamoxifen Series (UBC-TAM):  These cases were drawn from a
108    well-annotated cohort of patients treated with adjuvant tamoxifen without chemotherapy[6].  A
109    total of 625 of 632 (98.8%) patient samples that fully met study criteria passed a minimum
110    sequencing quality cutoff of at least 80% of targeted bases covered at greater than 20X (mean
111    coverage: 133X) with other quality metrics described in the supplementary data (Supplementary
112    Figure 1-5 and Supplementary Data 2). The final patient population had an average age of 67 at
113    diagnosis (range: 40-89+).  All were treated with five years of adjuvant tamoxifen, and were
114    primarily postmenopausal, grade 2 or 3 cancers, of ductal histologic subtype (Supplementary
115    Table 2). All were ER+ and at least 88.6% were clinically HER2- (13/625 unknown).  A subset of
116    463 of these patients had PAM50 subtyping data available from a previous study [6]. The median
117    follow up in the cohort examined was 25 years and one month.
118
119    POLAR cohort: This patient series was a case-control study of ER+ breast tumors, 175 of 194
120    (90.2%) patient samples passed minimum sequencing quality thresholds. A case was defined
121    as any patient who relapsed during follow-up, and controls were defined as lacking relapse
122    through a similar follow-up duration. Based on these definitions, there were 91 cases and 84
123    controls. Of the cases, 43 were early relapses (<5 years since diagnosis) and 48 were late
124    relapses (>5 years). Patients were only included if they received adjuvant endocrine therapy,
125    but chemotherapy was not an exclusion criterion, nor was menopausal status. These cases
126    were used in the mutation landscape and hotspot analyses only.
127
128    NCIC-MA12 Trial cohort.  These cases were drawn from a clinical trial in premenopausal
129    women treated with a standard adjuvant chemotherapy regimen and randomized to tamoxifen
130    versus observation. A total of 459 patient samples passed the minimum sequencing quality
131    threshold, of which 328 were hormone receptor positive (HR+), and only the HR+ cohort are
132    included here for most analyses. The majority were premenopausal (mean age of 45). All
133    patients received chemotherapy, and 48% were treated with 5 years of adjuvant tamoxifen. A
134    subset of 255 of these patients had PAM50 subtyping data available. The median follow up in
135    the cohort examined was 9.7 years
136
137    Across the three cohorts, there were 1,259 patient samples that passed minimum sequencing
138    quality thresholds and 1,128 of these were ER+ (UBC-TAM and POLAR) or HR+ (MA12).
139
140    Variant calling and filtering
141
142    A total of over 62 million variants were called in UBC-TAM. After extensive filtering against a set
143    of nearly 70,000 unmatched normal samples and manual review to eliminate common
144    polymorphisms and false positives (see methods), 1,991 putative somatic variants were
145    identified (0 to 26 variants per patient). A set of 1,693 mutations was defined as the "non-silent"
146    set for further analysis that excluded sequencing variants in splice regions, RNA genes (except
147    *MALAT1*), UTRs, introns, and all silent mutations. Finally, a set of 408 frameshift or nonsense
148    mutations was defined. The same filtering method was applied to both the POLAR and MA12
149    datasets. A total of 540 putative somatic mutations (436 non-silent, 145 FS/NS) were identified
150    in POLAR, and 2,104 (1,753 non-silent, 610 FS/NS) in MA12. Full details on these variants are
151    included in Supplementary Data 3 and summarized for key genes in Supplementary Figure 6.
152
153

2

154
155
156     Mutation landscape analysis.
157     In 1128 samples passing quality control standards, considering only non-silent mutations, 17
158     genes were mutated at a rate greater than 5%, and 6 at a rate greater than 10%; *PIK3CA* was
159     the only gene mutated in greater than 20% of samples (**Figure 1A**). The order from most
160     recurrent to least for the 10 most frequently mutated genes was: *PIK3CA* (41.1%), *TP53*
161     (15.5%), *MLL3* (13.4%), *MAP3K1* (12.0%), CDH1 (10.5%), MALAT1 (10.0%), GATA3 (9.1%),
162     MLL2 (8.7%), ARID1A (7.2%), and BRCA2 (6.6%). This list correlates well with previously
163     reported recurrently mutated genes. For example, the top 4 most significantly mutated genes in
164     the ER+ subset of TCGA breast project[3] were *PIK3CA* (24.3%), *TP53* (14.6%), *GATA3* (8.9%)
165     and *MAP3K1* (6.2%). The overall average mutation rate was estimated as 3.3 per MB of coding
166     sequence (range: 0.5 to 13.8 mutations per MB, excluding samples with no mutations called). In
167     order to determine whether mutations in any gene pair were mutually exclusive or co-occuring in
168     this dataset, a pairwise Chi-squared or Fisher's exact test was performed. Mutations in PIK3CA
169     and MAP3K1 were significantly more likely to co-occur (after BH FDR correction) in TAM
170     dataset, and were near significance in MA12 although not after correction (p = 0.08). These
171     results are summarized in Supplementary Data 4.
172
173     Hotspot analysis
174
175     As anticipated[7], mutations in *PIK3CA* at *E542K*, *E545K*, and *H1047R* were highly recurrent in
176     this study with 69/1259 (5.5%) E542K, 104 (8.3%) E545K, and 181 (14.4%) H1047R mutations
177     (Supplementary Figure 6C). Mutations in the ligand binding domain of *ESR1* (1.1%) were
178     extremely rare[3] (Supplementary Figure 6A). To uncover novel hotspots in these data, both Chi-
179     squared and Fisher's exact tests were performed using mutation frequencies from previous
180     sequencing studies as the expected values (see Methods for definition of multi-study MAF file)
181     (Supplementary Table 3). The most notable novel finding was in *CBFB* (**Figure 1B**). At least 6
182     different genomic alterations were observed in 15 patients (Supplementary Data 3) that affected
183     the donor splice site of exon 2. Manual review of this splice site identified at least two additional
184     patients with evidence for mutations at this location. The predicted effect of these mutations is
185     skipping of exon 2 or alternate donor site usage, each likely resulting in loss-of-function of the
186     *CBFB* protein. Additional splice site mutations were observed at the exon 2, exon 4 and exon 5
187     acceptor sites of *CBFB*. ErbB2 expressed the anticipated profile of activating mutations from
188     earlier publications[8] with 22/1259 (1.7%) samples harboring known activating mutations and
189     another 6 variants of unknown significance in the kinase domain or at the S310 residue (**Figure
190     8C**).
191
192     Somatic mutation association with PAM50-based intrinsic subtype
193
194     The PAM50 intrinsic subtype calls were obtained from previously published analyses to
195     compare their mutational profiles between UBC-TAM and the MA12 studies. In both studies
196     about half the patients had luminal A tumor. However, the MA12 cohort had a higher proportion
197     of non-luminal subtypes, with 19.8% HER2-E and 6.6% basal and fewer luminal B tumors
198     (25.1% versus 42.4%) (**Figure 2A-B**). Age density plots by subtype serve to emphasize the
199     large difference in the median age between the two sample cohorts (43 versus 65), and also the
200     influence of age with respect to the intrinsic subtype incidence. Namely, in the younger MA12
201     cohort, there is a younger peak incidence with basal-like breast cancer than Luminal A disease
202     (**Figure 2D**). In contrast in the older UBC-TAM cohort, an influence of age on intrinsic subtype
203     was not observed (**Figure 2C**). Relationships between intrinsic subtype and mutation patterns
204     were also explored, classifying mutation positive status as "non-silent", "missense",

3

205      nonsense/frame-shift (FS/NS) or FS/NS+splice site (Supplementary Data 5). The FDR
206      corrected p-value (q-value) took into account that 83 genes were examined. However, this level
207      of false discovery detection could be viewed as overly conservative in an exploratory analysis
208      and any gene mutation with q-value association of <0.2 was therefore considered reportable [9-
209      11]. For MA12, non-silent TP53 mutation was highly subtype-associated because of the very
210      high incidence in non-luminal versus luminal subtypes. PIK3CA and MAP3K1 mutations were
211      associated with Luminal A disease in both cohorts (Supplementary Figure 7A). Finally, there
212      was a strong association between Luminal B status and non-silent (Supplementary Figure 7B)
213      as well as FS/NS mutations in GATA3 (Supplementary Data 5, q value = 0.006). GATA3
214      mutations were present in 28-30% of Luminal B cases and less so in luminal A cases (5%).
215      Considering q values of <0.2 the associations between FS/NS and non-silent mutations in ATM
216      and Luminal B tumors in MA12 (8-13%) suggests that ATM loss is also a possible luminal B
217      driver (Supplementary Figure 7B), at least in younger women (MA12). Relationships between
218      age and mutation incidence were therefore also explored (Supplementary Figure 7C), with the
219      finding that both ATM mutation and GATA3 mutations were associated with an earlier age of
220      onset within the luminal B category (**Figure 2E and 2F**). Finally, NF1 mutations were
221      associated with the HER2-enriched subtype in the UBC-Tam series, explaining the association
222      with poor outcomes (Supplementary Figure 7B).
223
224      Survival analysis according to somatic mutation.
225
226      For the UBC-TAM Series (**Figure 3A**), univariate analysis of favorable prognostic associations
227      for breast cancer specific survival (BCSS) were detected for non-silent mutations in *MAP3K1*,
228      *ERBB3*, XBP1 and PIK3CA (**Figure 3B**, Supplementary Data 6). Adverse prognostic effects
229      were observed for non-silent mutations in *DDR1* and *TP53*, as well as for frame-shift and
230      nonsense (FS/NS) mutations in NF1. An analysis for recurrence free survival (RFS) produced
231      similar results, except for ARID1B, which was marginally associated with more favorable
232      outcome. A multivariate Cox model was applied to put each gene in the context of clinical
233      parameters (grade, tumor size and node status). These analyses indicated that the prognostic
234      effects of non-silent DDR1, PIK3CA, GATA3 FS/NS, TP53 and MAP3K1 mutations were
235      independent of grade and pathological stage (**Figure 3C**). Multiple correction testing, yielded
236      DDR1 as the only gene that remained significant with a q-value of 0.0003. (Supplementary
237      Data 5). For the MA12 clinical trial cohort (**Figure 4A**) we focused on overall survival
238      associations as this was the primary endpoint of the study and the most robust endpoint. A
239      number of rarely mutated genes were associated with poor outcome in univariate analysis as
240      displayed in **Figure 4B**. Multiple testing corrections indicated none of these findings could be
241      considered significant [9-11]. However, in multivariate analysis, based on the uncorrected p value,
242      the prognostic effects of mutations in ErbB2, ErbB4, LTK FS/NS, MAP3K4, PIK3R1, RB1, RELN
243      and TGFB2 were independent of pathological stage and grade (**Figure 4B**).
244
245      Verification of Prognostic effects of Mutations in METABRIC data.
246
247      While few genes were significant in univariate analysis after multiple testing correction, they
248      provide valuable hypotheses for further testing and validation. We therefore sought additional
249      data in the public domain to further assess the uncorrected p value-based findings in our data
250      set. The METABRIC consortium have reported somatic mutations in cBioPortal [12] with co-
251      reported detailed hormone receptor status, age at diagnosis (median age=64 years for ER+
252      patients), mean follow up of >8 years and disease-specific outcome [13, 14]. This data set provided
253      the opportunity to conduct a validation exercise for overlapping genes in the two data sets. For
254      the UBC-TAM series (**Figure 3**), 9 genes with a univariate p value of <0.05 were brought
255      forward for validation (**Figure 5**). Of the 6 overlapping genes also examined in METABRIC,

4

256 consistent prognostic effects independent of clinical variables were observed for non-silent
257 mutations in three genes, *MAP3K1* (favorable), *TP53* (unfavorable) and *NF1* FS/NS mutations
258 (unfavorable). For the MA12 series (**Figure 4**), 5 shared genes were identified with univariate p
259 values of <0.05, yet only *PIK3R1* mutations (non-silent or FS/NS) showed consistent adverse
260 prognostic effects (**Figure 6**). The Kaplan Meier survival plots for the consistent adverse
261 prognostic effects of *NF1* FS/NS and non-silent *PIK3R1* mutations are illustrated in **Figure 7A-**
262 **D**.
263
264 Prognostic interactions between PIK3CA and MAP3K1.
265
266 Since PIK3CA and MAP3K1 mutations co-associate, the combined effect of non-silent
267 mutations in these genes was examined. Patients with tumors exhibiting both genes mutated
268 have a more favorable clinical course than either singly mutant cases or cases without either
269 gene mutated. While the prognostic effects were strongest in the UBC-TAM series, this result
270 was also reproduced in the METABRIC data (**Figure 7E-F**).
271
272 Mutation Analyses for Uncommon Targetable Kinases.
273
274 Of the 83 genes analyzed, at least 8 are directly targetable with small molecules or antibodies
275 that are either FDA approved or in late-stage development (**Figure 8**). Pre-existing data on
276 these mutations is summarized (Supplementary Data 7). PIK3CA is not further discussed here,
277 since the mutation spectrum is well-described and large therapeutic studies are already
278 underway. A total of 23 patients with breast cancer with ErbB2 activating mutations were
279 identified. An examination of their locations revealed that ErbB2 mutations were, as expected,
280 clustered in 2 major domains, with 2 of 23 having extracellular domain mutations at residue 310
281 and 21 of 23 having kinase domain mutations between residues 755-842 [8, 15]. To further
282 investigate the preliminary finding of an adverse prognostic effect for ErbB2 mutation in the
283 MA12 series, an examination of the METABRIC data indicated that known activating mutations
284 in ErbB2 were associated with a near significant adverse effect (HR=1.71, P=0.075)
285 (Supplementary Figure 8). For ERBB3, 2 known-activating mutations were identified (V104L
286 and E928A)[16]. The DDR1 kinase domain mutation, R776W, is possibly homologous to EGFR
287 hot spot mutation L858R, but the remaining DDR1 variants are of unknown significance. For the
288 mutations in JAK1, 3 of 12 are loss of function mutations (frame shift or non-sense) and the
289 S816* mutation has been reported in a lung adenocarcinoma sequencing data set [17]. The loss
290 of function mutations in JAK1 have been shown to associate with immune therapy resistance [18,
291 19]. A few mutations identified in ERBB4, MET, and PDGFRA have been previously reported but
292 those reported here have not been functionally tested.
293
294
295 Discussion
296
297 The landscape of recurrently mutated genes in ER+ breast cancer observed in this study is
298 consistent with reports where matched germline samples were available, indicating that our
299 variant filters were effective for somatic mutation detection in a research setting. Overall,
300 mutation rates were higher in our cohort (e.g., for *PIK3CA*, *MLL3*, *MAP3K1*) than the TCGA
301 cohort, but were also lower for a few specific genes (e.g., *TP53* and *GATA3*). Due to higher
302 sequencing data coverage of recurrently mutated target genes than TCGA and the use of a
303 different hybrid capture reagent, we were likely able to detect mutations that were missed with
304 lower-depth exome or whole genome sequencing data. It is also possible that in some instances
305 we overestimated somatic mutation rates, due to the lack of matched normal samples and
306 imperfections in our germline polymorphism filtering. In particular, a significant number of

5

307  BRCA1 and BRCA2 mutations are likely *de novo* germline mutations that we would not be able
308  to easily distinguish from somatic mutations. Of the 117 non-silent BRCA1/2 mutations
309  observed (from 110/1128 patients; 7 patients had two hits) 74 were observed at a VAF greater
310  than 40% and 31 were greater 60%. Variants with VAFs this high are less likely to be somatic
311  given the general expectation of impure tumor samples and heterozygous mutations. Indeed,
312  the VAFs for BRCA1/2 non-silent mutations (mean=44.8%) were significantly higher than for
313  other genes (mean=36.4%, p = 6.96e-06). There were 8 known pathogenic (ENIGMA expert
314  reviewed) mutations according to a search of the BRCA Exchange database
315  (http://brcaexchange.org, Nov 12, 2017) and another 37 likely pathogenic (FS/NS) mutations. Of
316  the remaining, 4 were known benign according to expert review (ENIGMA), and 8 benign, 15
317  likely benign and 45 variants of unknown significance according to all public sources.
318
319  The discovery of a novel recurrent *CBFB* (core binding factor subunit beta) splice site mutation
320  in this cohort illustrates a limitation of exome capture reagents. The affected bases in exon 2 of
321  CBFB display reduced sequence coverage, possibly due to high GC content, in the breast
322  TCGA exome dataset (Supplementary Figures 9-10). This site was mutated in at least 1.5% of
323  ER+ breast cancers sequenced, bringing the overall rate of CBFB mutations to nearly 6%,
324  which should drive further investigation of this gene in ER+ breast cancer pathogenesis. *CBFB*
325  functions as a subunit in a heterodimeric core binding transcription factor that interacts with
326  *RUNX1*[20]. Consistent with this model, *CBFB* mutants were mutually exclusive from *RUNX1*
327  mutants in this cohort with only a single sample harboring non-silent mutations in both *CBFB*
328  and *RUNX1*.
329
330  The UBC-TAM and MA12 studies revealed different lists of potentially prognostic mutations.
331  Prognostic effects are likely to be strongly affected by the use of systemic therapy as well as by
332  patient age at diagnosis.  The UBC-TAM series is the simplest study to interpret from a drug
333  resistance perspective since the only systemic therapy was tamoxifen.  Thus, the consistent
334  adverse effect of NF1 FS/NS mutation on prognosis is intriguing as this result is consistent with
335  results from an *in vitro* screen for tamoxifen resistance[21].  Understanding why only FS/NS
336  mutations predict poor outcome, rather than missense or other non-silent mutations, will require
337  further investigation.  In contrast, PIK3R1 mutation emerged as a consistent poor prognosis
338  mutation from the MA12 analysis, with validation in METABRIC.  The proposed favorable
339  prognostic effects of PIK3CA mutation were observed in the UBC-TAM series, but were not
340  found to be independent of stage and grade, and PTEN mutations were neutral.
341
342  According to our validation results, NF1, PIK3R1, MAP3K1, PIK3CA and TP53 are likely to be
343  prognostic drivers. In postmenopausal women treated with adjuvant endocrine therapy, DDR1,
344  PRKDC and XBP1 should be further studied and of these DDR1 is the strongest candidate
345  because it was significant despite strict false discovery correction. DDR1 is a collagen-binding
346  receptor expressed in epithelial cells that stabilizes E-cadherin–mediated intracellular
347  adhesion[22]. *DDR1* mutations also occur in endometrial cancer[23], acute leukemia[24] and lung
348  cancer[25]. Loss of DDR1 (DDR1-null mice) produces hyper-proliferation and abnormal branching
349  of mammary ducts, suggesting DDR1 is a breast tumor suppressor[26]. The relationship between
350  truncating mutations in NF1 and poor outcome is consistent with an siRNA screen for genes
351  whose loss generates tamoxifen resistance[21].  Mutations in PRKDC will potentially produce a
352  defective ATM response/low ATM levels [27] which is interesting in the context of the finding
353  herein that ATM mutations are a potential luminal B driver gene.  The significance of a defective
354  ATM pathway as a cause of endocrine resistance is highlighted by the recent finding that
355  dysregulation of the MutL complex (MLH1, PMS1 and PMS2) causes failure of ATM/CHK2-
356  based negative regulation of CDK4/6 [28]. Prognostic candidate mutations revealed by the MA12
357  analysis were different from the UBC TAM series, likely reflecting the different patient profiles

358 and adjuvant treatments illustrated in **Figure 2**. The prognostic effects of mutations ERBB2,
359 ERBB4, JAK1, LTK, MAP3K4, MET, PDGFRA, RB1, RELN, TGFB2, all await further study with
360 even larger sample sizes.
361
362 In conclusion, we have successfully utilized clinically well-annotated, uniformly treated patient
363 samples using DNA from archival material greater than 20 years old that lacks a matched
364 normal to explore the prognostic effects encoded by the mutational landscape of ER+ breast
365 cancer. We were able to confirm our prospective hypothesis that MAP3K1 is associated with
366 indolent disease and TP53 with adverse outcomes. We also associated NF1 FS/NS mutations
367 with strong adverse effects on prognosis. Similarly, PIK3R1 mutations were associated with an
368 adverse prognosis in contrast to PIK3CA mutation. This suggests somatic mutations in these
369 two physically interacting gene products are not biologically equivalent with respect to PI3
370 kinase pathway activation and resistance effects. The possibility that the long tail of low
371 frequency mutation events in luminal type breast cancer may harbor multiple molecular
372 explanations for poor outcomes is an important finding that should spur collaborative efforts to
373 thoroughly screen thousands of properly annotated cases. Only after these iterative efforts of
374 proposing and confirming candidates will a clinically useful and comprehensive somatic
375 mutation-based classification of ER+ breast cancer emerge. In the meantime, functional studies
376 should be pursued to understand the biological effects of somatic mutations, prioritizing these
377 studies according to whether the mutations are driving an adverse prognostic effect.
378
379 Methods
380
381 For the UBC-TAM series, an institutional review board approved study was based on formalin-
382 fixed paraffin embedded (FFPE) primary tumor blocks from 947 female patients diagnosed with
383 estrogen receptor positive invasive breast cancer in the province of British Columbia in Canada
384 between 1986 and 1992[6, 29-31]. The sample flow and analysis are provided in a REMARK
385 summary (**Figure 3A**). DNA was isolated from tumor-rich regions using the Qiagen blood and
386 tissue kit, which yielded sufficient DNA in 645 samples, of which 625 met all study criteria and
387 had sufficient sequence coverage. Similarly, approved studies provided 194 and 454 HR+
388 patient samples for the POLAR and MA12 (**Figure 4A**) cohorts. A total of 175 POLAR and 459
389 (328 HR+) MA12 samples yielded sufficient DNA and had sufficient sequence coverage for
390 analysis. Detailed descriptions of the patient data sets are provided in Supplementary Table 3.
391 A meta-analysis of six existing published large-scale breast cancer sequencing studies [1-3, 5, 32, 33]
392 was performed to identify genes with recurrent coding region somatic mutations in breast cancer
393 (Supplementary Data 1). Additional drug targets[34] and genes with relevance to breast cancer
394 from targeted sequencing[35], copy-number studies[13] or knowledge relating to somatic or germline
395 mutations (e.g., *BRCA1*, *BRCA2*, *ERBB2*, *ESR1* and *PRLR*) were also included. This resulted in
396 a final list of 83 breast-cancer-related genes (Supplementary Table 1). These genes were
397 targeted comprehensively with 3,029 complementary probes for hybridization-based enrichment
398 (Supplementary Data 8). Sequencing libraries were constructed, hybridized with capture probes,
399 multiplexed and run on a single flow cell with up to 96 samples per pool per lane yielding
400 approximately 375 Mb of DNA sequence per sample from an Illumina HiSeq paired end 2 X
401 100bp (TAM) or 2 X 125bp (POLAR, MA12) sequencing run following manufacturer's protocols.
402
403 Variant calling was performed with the Genome Modeling System as previously described[36].
404 Specifically, sequence data were aligned to reference sequence build GRCh37 using BWA[37]
405 and de-duplicated with Picard. SNVs and indels were detected using the union of samtools[38]
406 and VarScan2[5] and annotated using Ensembl version 70. Variants were restricted to the coding
407 regions of targeted genes and filtered for false positives and germline polymorphisms against a
408 database of nearly 70,000 unmatched normals from the ExAC consortium, 1000 Genomes[39,]

7

409    NHLBI exomes[40] and TCGA data sets[3, 41]. A binomial probability model was then applied to the
410    variants using VAF and total coverage to determine a log-likelihood ratio of being a somatic
411    variant as previously described[42] (See Supplementary Methods). After filtering, all remaining
412    variants were manually reviewed. To ensure that variants of known clinical relevance were not
413    missed by automated variant calling approaches, a knowledge-based variant calling strategy
414    was performed focused on the mutations in the Database of Curated Mutations[43].
415
416    Patient groups were defined by mutation status or truncating mutation status for each gene.
417    Fisher's exact and Chi-squared tests were used for hotspot analysis, mutual exclusivity or co-
418    occurrence, and other categorical clinical statistics (e.g., mutation status vs. intrinsic subtype) as
419    appropriate. Univariate Kaplan-Meier and Cox survival analyses were performed for breast-
420    cancer-specific survival (BCSS), relapse free survival (RFS), or overall survival (OS) with non-
421    silent or truncating mutation status as a factor. Significant survival differences between the
422    groups were determined by log rank (Mantel-Cox) test. The Benjamini-Hochberg method was
423    performed for multiple testing corrections to report the false discovery rate adjusted p-value (q-
424    value). A multivariate Cox proportional hazard model was fitted to BCSS and RFS separately on
425    gene mutation status, node status, grade and tumor size and adjusted hazard ratios were
426    calculated with Wald test p-values. All statistical analyses were performed in the R statistical
427    programming language with core, 'survival' and 'multtest' libraries. Genomic visualizations were
428    created with ProteinPaint[44] and GenVisR[45].
429
430
431
432
433

**References**

434
435
436    1.    Banerji, S. et al. Sequence analysis of mutations and translocations across breast
437          cancer subtypes. *Nature* **486**, 405-409 (2012).
438    2.    Stephens, P.J. et al. The landscape of cancer genes and mutational processes in breast
439          cancer. *Nature* **486**, 400-404 (2012).
440    3.    Cancer Genome Atlas, N. Comprehensive molecular portraits of human breast tumours.
441          *Nature* **490**, 61-70 (2012).
442    4.    Kennecke, H.F. et al. Late risk of relapse and mortality among postmenopausal women
443          with estrogen responsive early breast cancer after 5 years of tamoxifen. *Ann Oncol* **18**,
444          45-51 (2007).
445    5.    Ellis, M.J. et al. Whole-genome analysis informs breast cancer response to aromatase
446          inhibition. *Nature* **486**, 353-360 (2012).
447    6.    Nielsen, T.O. et al. A comparison of PAM50 intrinsic subtyping with
448          immunohistochemistry and clinical prognostic factors in tamoxifen-treated estrogen
449          receptor-positive breast cancer. *Clin Cancer Res* **16**, 5222-5232 (2010).
450    7.    Samuels, Y. et al. High frequency of mutations of the PIK3CA gene in human cancers.
451          *Science* **304**, 554 (2004).
452    8.    Bose, R. et al. Activating HER2 mutations in HER2 gene amplification negative breast
453          cancer. *Cancer Discov* **3**, 224-237 (2013).
454    9.    Amar, D., Shamir, R. & Yekutieli, D. Extracting replicable associations across multiple
455          studies: Empirical Bayes algorithms for controlling the false discovery rate. *PLoS*
456          *Computational Biology* **13**, e1005700 (2017).
457    10.    Capanu, M. & Seshan, V.E. False discovery rates for rare variants from sequenced data.
458          *Genetic epidemiology* **39**, 65-76 (2015).
459    11.    Efron, B. Size, Power and False Discovery Rates. *The Annals of Statistics* **35**, 1351-
460          1377 (2007).
461    12.    Gao, J. et al. Integrative analysis of complex cancer genomics and clinical profiles using
462          the cBioPortal. *Sci Signal* **6**, pl1 (2013).
463    13.    Curtis, C. et al. The genomic and transcriptomic architecture of 2,000 breast tumours
464          reveals novel subgroups. *Nature* **486**, 346-352 (2012).
465    14.    Pereira, B. et al. The somatic mutation profiles of 2,433 breast cancers refines their
466          genomic and transcriptomic landscapes. *Nat Commun* **7**, 11479 (2016).
467    15.    Ma, C.X. et al. Neratinib Efficacy and Circulating Tumor DNA Detection of HER2
468          Mutations in HER2 Nonamplified Metastatic Breast Cancer. *Clin Cancer Res* **23**, 5687-
469          5695 (2017).
470    16.    Jaiswal, B.S. et al. Oncogenic ERBB3 mutations in human cancers. *Cancer Cell* **23**,
471          603-617 (2013).
472    17.    Imielinski, M. et al. Mapping the hallmarks of lung adenocarcinoma with massively
473          parallel sequencing. *Cell* **150**, 1107-1120 (2012).
474    18.    Shin, D.S. et al. Primary Resistance to PD-1 Blockade Mediated by JAK1/2 Mutations.
475          *Cancer Discov* **7**, 188-201 (2017).
476    19.    Zaretsky, J.M. et al. Mutations Associated with Acquired Resistance to PD-1 Blockade in
477          Melanoma. *N Engl J Med* **375**, 819-829 (2016).
478    20.    Lukasik, S.M. et al. Altered affinity of CBF beta-SMMHC for Runx1 explains its role in
479          leukemogenesis. *Nat Struct Biol* **9**, 674-679 (2002).
480    21.    Mendes-Pereira, A.M. et al. Genome-wide functional screen identifies a compendium of
481          genes affecting sensitivity to tamoxifen. *Proc Natl Acad Sci U S A* **109**, 2730-2735
482          (2012).

483   22.   Yeh, Y.C., Wu, C.C., Wang, Y.K. & Tang, M.J. DDR1 triggers epithelial cell
484         differentiation by promoting cell adhesion through stabilization of E-cadherin. *Molecular*
485         *biology of the cell* **22**, 940-953 (2011).
486   23.   Rudd, M.L. et al. Mutational analysis of the tyrosine kinome in serous and clear cell
487         endometrial cancer uncovers rare somatic mutations in TNK2 and DDR1. *BMC Cancer*
488         **14**, 884 (2014).
489   24.   Loriaux, M.M. et al. High-throughput sequence analysis of the tyrosine kinome in acute
490         myeloid leukemia. *Blood* **111**, 4788-4796 (2008).
491   25.   Ding, L. et al. Somatic mutations affect key pathways in lung adenocarcinoma. *Nature*
492         **455**, 1069-1075 (2008).
493   26.   Vogel, W.F., Aszodi, A., Alves, F. & Pawson, T. Discoidin domain receptor 1 tyrosine
494         kinase has an essential role in mammary gland development. *Molecular and cellular*
495         *biology* **21**, 2906-2917 (2001).
496   27.   Peng, Y. et al. Deficiency in the catalytic subunit of DNA-dependent protein kinase
497         causes down-regulation of ATM. *Cancer Res* **65**, 1670-1677 (2005).
498   28.   Haricharan, S. et al. Loss of MutL Disrupts CHK2-Dependent Cell-Cycle Control through
499         CDK4/6 to Promote Intrinsic Endocrine Therapy Resistance in Primary Breast Cancer.
500         *Cancer Discov* **7**, 1168-1183 (2017).
501   29.   Cheang, M.C. et al. Ki67 index, HER2 status, and prognosis of patients with luminal B
502         breast cancer. *J Natl Cancer Inst* **101**, 736-750 (2009).
503   30.   Liu, S. et al. Prognostic significance of FOXP3+ tumor-infiltrating lymphocytes in breast
504         cancer depends on estrogen receptor and human epidermal growth factor receptor-2
505         expression status and concurrent cytotoxic T-cell infiltration. *Breast Cancer Res* **16**, 432
506         (2014).
507   31.   Parker, J.S. et al. Supervised risk predictor of breast cancer based on intrinsic subtypes.
508         *J Clin Oncol* **27**, 1160-1167 (2009).
509   32.   Kan, Z. et al. Diverse somatic mutation patterns and pathway alterations in human
510         cancers. *Nature* **466**, 869-873 (2010).
511   33.   Shah, S.P. et al. The clonal and mutational evolution spectrum of primary triple-negative
512         breast cancers. *Nature* **486**, 395-399 (2012).
513   34.   Griffith, M. et al. DGIdb: mining the druggable genome. *Nat Methods* **10**, 1209-1210
514         (2013).
515   35.   Chanock, S.J. et al. Somatic sequence alterations in twenty-one genes selected by
516         expression profile analysis of breast carcinomas. *Breast Cancer Res* **9**, R5 (2007).
517   36.   Griffith, M. et al. Genome Modeling System: A Knowledge Management Platform for
518         Genomics. *PLoS Comput Biol* **11**, e1004274 (2015).
519   37.   Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler
520         transform. *Bioinformatics* **25**, 1754-1760 (2009).
521   38.   Li, H. et al. The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**,
522         2078-2079 (2009).
523   39.   Genomes Project, C. et al. A map of human genome variation from population-scale
524         sequencing. *Nature* **467**, 1061-1073 (2010).
525   40.   Fu, W. et al. Analysis of 6,515 exomes reveals the recent origin of most human protein-
526         coding variants. *Nature* **493**, 216-220 (2013).
527   41.   Cancer Genome Atlas Research, N. Genomic and epigenomic landscapes of adult de
528         novo acute myeloid leukemia. *N Engl J Med* **368**, 2059-2074 (2013).
529   42.   Krysiak, K. et al. A genomic analysis of Philadelphia chromosome-negative AML arising
530         in patients with CML. *Blood Cancer J* **6**, e413 (2016).
531   43.   Ainscough, B.J. et al. DoCM: a database of curated mutations in cancer. *Nat Methods*
532         **13**, 806-807 (2016).

533   44.   Zhou, X. et al. Exploring genomic alteration in pediatric cancer using ProteinPaint. *Nat*
534         *Genet* **48**, 4-6 (2016).
535   45.   Skidmore, Z.L. et al. GenVisR: Genomic Visualizations in R. *Bioinformatics* **32**, 3012-
536         3014 (2016).
537
538
539
540
541
552
553
554   Contributions:
555
556   O.L.G., N.C.S., T.O.N., M.J.E., E.R.M. designed the experiments; M.G., J. K., C.A.M., K.K.,
557   J.H., B.J.A., Z.L.S., K.C. R.K. C.F., L.C., J.E.S., S.D., V.M., D.E.L., R.S.F., S.L., R.K.W.
558   generated the sequencing data, T.O.N., B.Y., M.D. S.L., and D.V. orchestrated the sample
559   pipeline., M.A., O.L.G. and N.C.S., prepared the figures and tables. M.A., J.L., and D.T.
560   provided statistical analysis.  S.M.K., R.B., and E.C.C. provided functional annotations., T.O.N
561   provided pathology analysis. M.J.E., N.C.S., M.A., and O.L.G. wrote the manuscript. E.R.M.,
562   T.O.N., M.D., critically read and commented on the manuscript.
563
564   Conflict of Interest:
565
566   Dr. Ellis and Dr. Mardis report income on patents on the PAM50 intrinsic subtype algorithm. Dr.
567   Ellis reports ownership in Bioclassifier LLC that licenses PAM50 patents to Nanostring for the
568   Prosigna breast cancer prognostic test. Commercial platforms and algorithms were not used in
569   the analyses reported in this paper.

570
571  Figure Legends
572
573  **Figure 1. Mutation recurrence and novel splice site mutation**
574  A) The overall mutation recurrence rate ranged from 41.1% of samples for *PIK3CA* to 0.0% for
575  *PIN1*. The figure depicts non-silent mutations for all 1128 patients for the top 16 most
576  recurrently mutated genes (>5% recurrence). If a patient had multiple mutations it is colored
577  according to the "most damaging" mutation following the order presented in the Mutation Type
578  legend (vertical color bar). Mutations per MB were calculated using the total number of
579  mutations observed over the total exome space corresponding to the tiled space from "SeqCap
580  EZ Human Exome Library v2.0". A correction factor was applied for genes not assayed using
581  the expected number of additional mutations based on TCGA data. B) Mutation recurrence
582  rates (amino acid level) in this study were compared to previously reported mutation rates from
583  a multi-study MAF file of six reported breast cancer sequencing studies (Supplementary Data
584  1). An entirely novel mutation "hot spot" was discovered affecting the exon 2 splice (donor) site
585  of *CBFB* in at least 15 patients. Six different single nucleotide substitutions, insertions and
586  deletions were observed, all affecting either the first or second base of the donor splice site.
587  These mutations were most likely missed in previous studies because of a lack of sequencing
588  coverage due to the GC-rich nature of exons 1 and 2 of *CBFB* (Supplementary Figures 9-10).
589  Such mutations are predicted to significantly alter the canonical donor site and result in either
590  alternate donor usage or skipping of one or more exons of *CBFB*.
591
592  **Figure 2. Cross-cohort age and subtype analysis**
593  A-B) Percentage composition of samples by intrinsic subtype of the tumor in the two discovery
594  cohorts for UBC-TAM (A) and MA12 (B) cohorts. C-D) Age-density plots for patients categorized
595  by intrinsic subtype in UBC-TAM (C) and MA12 (D) cohorts. The overall median age shows that
596  UBC-TAM is constituted mostly of post-menopausal patients (median age=65), in contrast to
597  MA12, which has younger patients (median age=43). E-F) Younger luminal B subtype patients
598  harbor GATA3 (E) and ATM (F) mutations in the combined set of UBC-TAM and MA12 Luminal
599  B cases (median age=52, p=0.01; median age=58, p=0.03 for GATA3 and ATM respectively).
600
601  **Figure 3. Candidate discovery from UBC-TAM cohort and prognosis evaluation**
602  (A) DNA was extracted from tumor specimens from 947 patients with ER+ breast cancer treated
603  with tamoxifen monotherapy for 5 years. 632 samples with adequate yield were sequenced for
604  83 genes known to be recurrently mutated or breast cancer relevant. A total of 625 samples
605  passed minimum quality checks and were sequenced to an average of 135.8X coverage. A total
606  of ~62 million variants from the reference genome were identified. Extensive filtering and
607  manual review reduced this list to 1,991 putatively somatic variants. Survival analysis was
608  applied to non-silent and truncating gene mutation status versus disease outcome (relapse or
609  breast-cancer-specific death). In addition, mutations were analyzed for novel hotspots, patterns
610  of mutual exclusivity or co-occurrence and association with clinical variables. (B) Forest plot of
611  impact of mutations in candidate genes, identified using UBC-TAM population, on breast-
612  cancer-specific-survival (red) and recurrence-free survival (blue). The variant types are
613  characterized based on non-silent or nonsense/frameshift (FS/NS) mutations. The box size is
614  relative to frequency of mutations listed in the analysis, larger boxes represent high incidence
615  rate mutations. (C) Multivariate forest plot of effect of mutations in UBC-TAM candidate genes
616  on breast cancer specific-survival when assessed together with clinical factors including Tumor
617  Grade, Node positivity and Tumor Size (>5cm).
618
619  **Figure 4. Candidate discovery from MA12 cohort and prognosis evaluation**

620    (A) DNA was extracted from tumor specimens and 470 samples with adequate yield were
621    sequenced for 83 genes known to be recurrently mutated or breast cancer relevant.  A total of
622    459 (328 HR+) samples passed minimum quality checks and were sequenced to an average of
623    272.6X coverage. A total of 406 million variants from the reference genome were identified.
624    Extensive filtering and manual review reduced this list to 2104 putatively somatic variants.
625    Survival analysis was applied to non-silent and truncating gene mutation status versus overall
626    survival. (B) Forest plot showing effect of mutation in candidate genes on overall survival
627    (univariate - blue, multivariate - orange), along with the clinical factors used in the multivariate
628    analysis, tumor grade, node positivity and tumor size (>5cm) in black. The box size is relative to
629    frequency of mutations listed in the analysis, larger boxes represents high incidence rate
630    mutations. Note: a few boxes are not shown if their hazard ratio were greater than 4.0.
631
632    **Figure 5. Validation of UBC-TAM candidates in ER+ METABRIC**
633    A) Six out of nine candidate genes from UBC-TAM analysis had mutations reported in
634    METABRIC cohort. 1060 ER+ samples with disease-specific survival information were used to
635    test the effect of mutations in the candidate genes on prognosis. B) Forest plot shows effect of
636    mutated candidate genes on disease-specific survival in METABRIC ER+ cohort with univariate
637    cox proportional-hazard ratio in blue and multivariate in orange.  The clinical factors used in the
638    multivariate analysis, namely tumor grade, node positivity and tumor size (>5cm), are shown in
639    black. The box size is relative to frequency of mutations listed in the analysis, larger boxes
640    represent genes with higher incidence rate of mutations.
641
642    **Figure 6. Validation of MA12 candidates in ER+ METABRIC**
643    A) Five out of eleven candidates from MA12 analysis had mutations reported in the METABRIC
644    cohort. 1415 ER+ samples with overall survival information was used to test the effect of
645    mutations in the candidate genes on prognosis. B) Forest plot shows effect of mutated
646    candidate genes, shortlisted based on MA12 mutation analysis, on overall survival in
647    METABRIC ER+ breast cancer patients. Univariate (blue) and multivariate (orange) cox
648    proportional-hazard ratio depict the independent prediction of survival outcomes for the six
649    candidate genes. The box size is relative to frequency of mutations listed in the analysis, larger
650    boxes represent genes with higher incidence rate of mutations.
651
652    **Figure 7. Kaplan-Meier plots**
653    A-B) Kaplan-Meier graph showing the prognostic role of NF1 mutations, separated by variant
654    type – Missense (MUT MS, green), Frameshift/Nonsense (MUT FS/NS, blue) in ER+ breast
655    cancer patients from A) UBC-TAM and B) METABRIC cohort establishing the association
656    between FS/NS mutations in NF1 with poor prognosis. C-D) Kaplan-Meier graph showing the
657    prognostic role of PIK3R1 in C) MA12 and D) METABRIC ER+ breast cancer patients,
658    categorized based on tumors with wildtype (WT, black) or mutated PIK3R1 non-silent mutations
659    (MUT, red). E-F) Kaplan-Meier graph demonstrating co-occurrence of non-silent mutations in
660    MAP3K1 and PIK3CA (red) in E) UBC-TAM and F) METABRIC associates with better survival
661    when compared against tumors with mutations exclusively in MAP3K1 (blue) or PIK3CA (green)
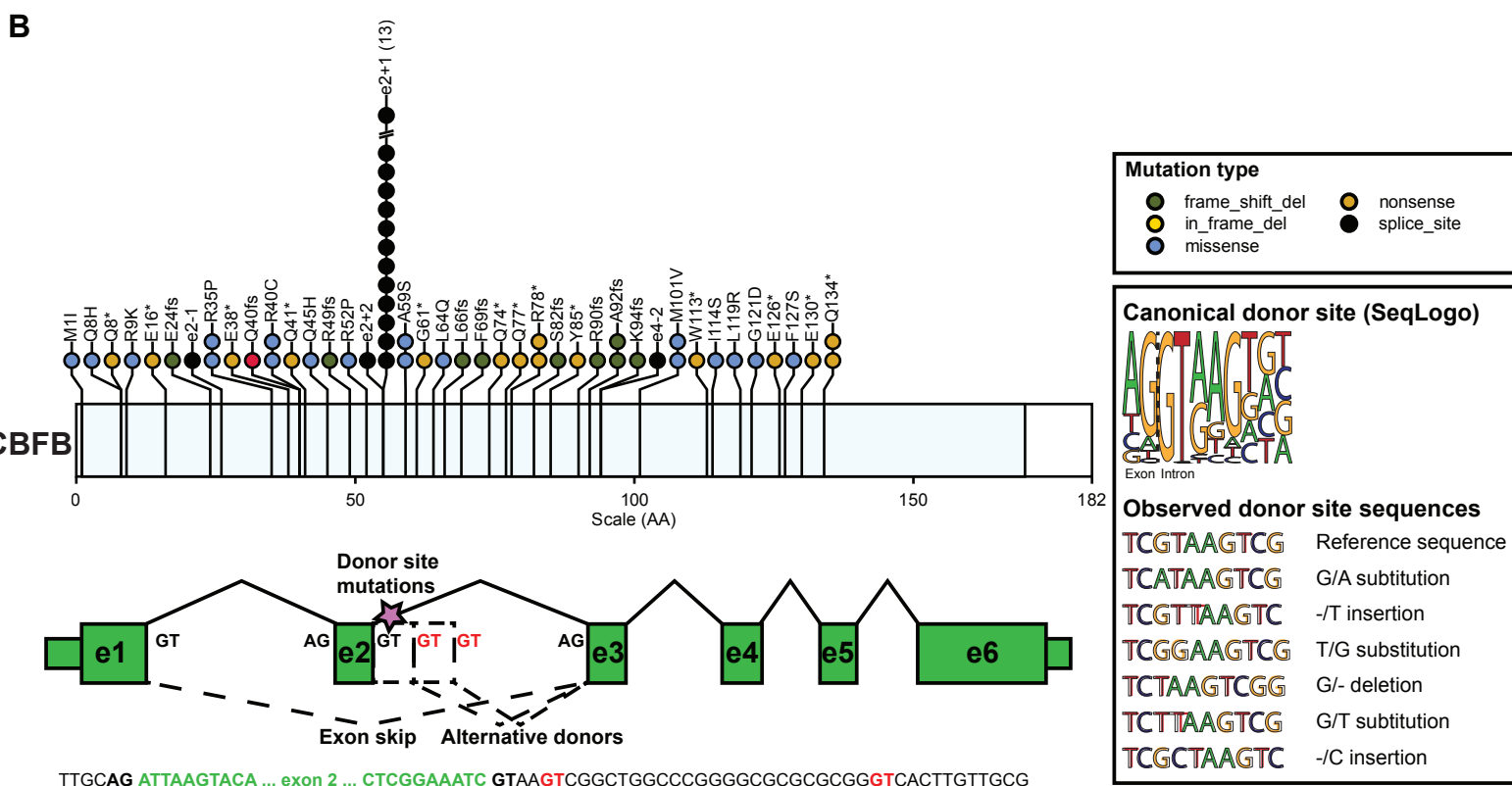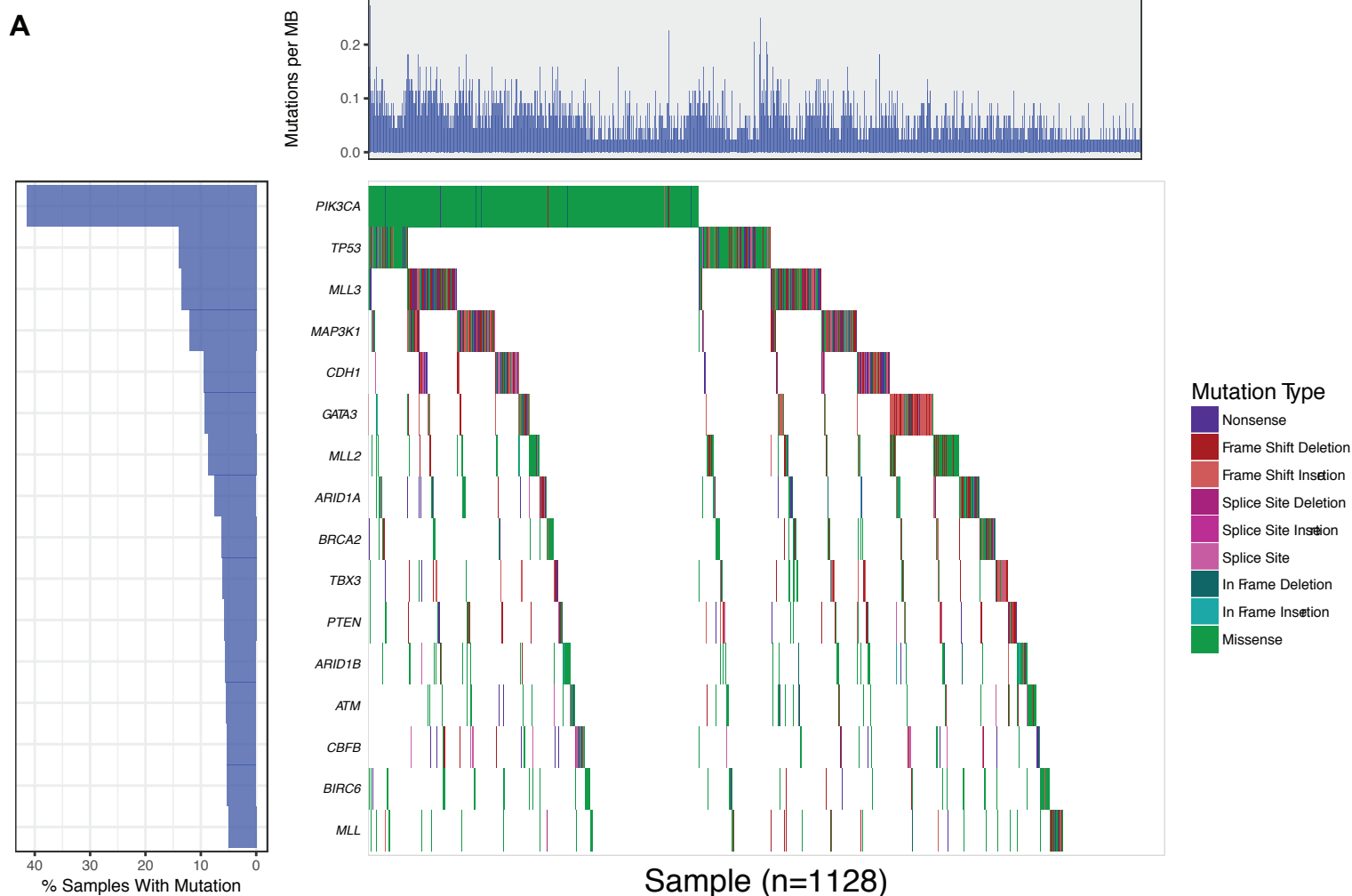662    or wildtype for both MAP3K1 and PIK3CA (black). p, log rank (Mantel-Cox) test p-value.
663
664    **Figure 8. Mutation profiles for selected genes**
665    Mutation frequency plots illustrate all non-silent mutations (TAM, POLAR, and MA12; n=1259)
666    for representative transcripts for several kinase genes of interest. The domains belonging to A)
667    DDR1 (RefSeq ID: NM_013994) and B) JAK1 (NM_002227) are indicated below the schematic
668    diagram of each gene. The ECD (extracellular domain), TM (transmembrane domain), and
669    kinase domain are depicted as green, red, and orange bars respectively for C) ERBB2
670    (NM_004448), D) ERBB3 (NM_001982), E) ERBB4 (NM_005235), F) MET (NM_000245), and

13

671    G) PDGFRA (NM_006206). The variant counts across the three datasets for each gene are
672    provided below the gene's name. Note, in the mapping from Ensembl (**Supplementary Data 3**)
673    to RefSeq annotations (required for use of ProteinPaint tool) a small number of variants
674    annotations may have changed or been lost, despite selecting the most similar representative
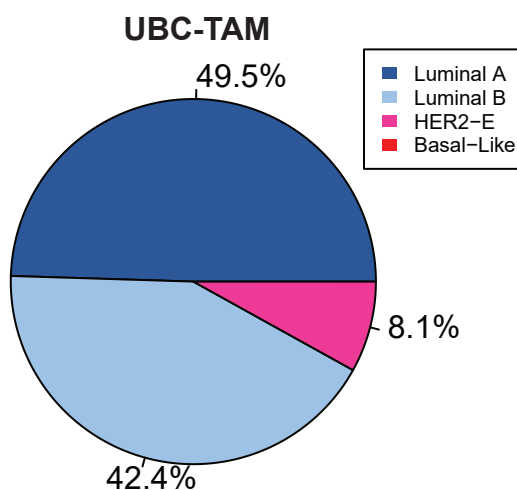675    transcript possible.
676

# Figure 1.

**A**

**B**

Figure 2.

**a**

**UBC-TAM**

49.5%

| Luminal A |
| Luminal B |
| HER2−E |
| Basal−Like |

8.1%

42.4%

**b**

**MA12**

48.6%

6.6%

25.1%

19.8%

**c**

Density

Her2-E

LumB

Overall
(Median=65)

LumA

Patient age in years

N = 529

**d**

Density

LumA

LumB

Overall
(Median=43)

Her2-E

Basal-Like

Patient age in years

N = 328

**e**

p=0.01

**GATA3 Mut**
(Median=52,
n=33)

**GATA3 WT**
(Median=66,
n=237)

Density

Age of patient

N=270 (LuminalB TAM+MA12)

**f**

p=0.03

**ATM Mut**
(Median=58,
n=20)

**ATM WT**
(Median=66,
n=250)

Density

Age of Patient

N=270 (LuminalB TAM+MA12)

# Figure 3.

**a**

**Patient**
- 947 Tam-treated patients.
- 645 with >50ng DNA.
- 632 ER+, tumor-only.
Gene Selection:
- 83 genes.
- 8 studies' recurrently mutated genes

- 625 samples meet coverage criteria.
- >60 million raw variant calls.
- Extensive filtering to remove germline calls.
- 1991 variants called as likely somatic.

- Mutation Landscape.
- Survival Analysis.
- Mutational Analysis.

**Data Analysis**

**Filtering Workflow**

Remove all variants with > .1% GMAF in 1000 genomes, NHLBI, ExAC. → Remove all artifacts seen using pipeline on 1063 exome and 87 WGS normal. → Knowledge-based variant detection using the DoCM database. → Manual review of all remaining variant calls.

**b**



| Gene | Variant | p-value |
|------|---------|---------|
| ARID1B | nonsilent | 0.1410 / 0.0310 |
| DDR1 | nonsilent | 2.42E-05 / 0.0047 |
| ERBB3 | nonsilent | 0.0606 / 0.0841 |
| MAP3K1 | nonsilent | 0.0033 / 0.0245 |
| NF1 | FS/NS | 0.1256 / 0.0167 |
| PIK3CA | nonsilent | 0.0260 / 0.0349 |
| PRKDC | nonsilent | 0.0380 / 0.0179 |
| TP53 | FS/NS | 0.0288 / 0.0127 |
| TP53 | nonsilent | 0.0216 / 0.0899 |
| XBP1 | nonsilent | 0.0683 / 0.0593 |

BCSS    RFS

**Univariate Hazard Ratio**

**c**



| Gene | Variant | p-value |
|------|---------|---------|
| ARID1B | nonsilent | 0.3145 |
| DDR1 | nonsilent | 0.0000 |
| ERBB3 | nonsilent | 0.0687 |
| GATA3 | FS/NS | 0.0364 |
| MAP3K1 | nonsilent | 0.0014 |
| MET | nonsilent | 0.0859 |
| NF1 | FS/NS | 0.2322 |
| PIK3CA | nonsilent | 0.0348 |
| PRKDC | nonsilent | 0.4917 |
| TP53 | FS/NS | 0.0333 |
| TP53 | nonsilent | 0.0389 |
| Tumor Grade | clinical | 0.0026 |
| Node positivity | clinical | 0.0000 |
| Tumor Size >5cm | clinical | 0.0000 |

**Multivariate Hazard Ratio**

**Figure 4.**

**Figure 5.**

**Figure 6.**



**a**

MA12 Candidates
*(UVA p<0.05)*

ERBB2    MAP3K4
ERBB4    MET
PIK3R1    PDGFRA
RB1    LTK
JAK1    TGFB2
RELN

METABRIC dataset

- Mutations extracted from cbioportal (Pereira, Nat Comm 2016)
- Clinical annotations linked using Oncomine Curtis Breast 2 dataset (Curtis, Nature 2012)

Univariate analysis MVA with Grade, Tumor Size, Node Status

Validation

Metabric Cohort (2369)

ER+ with OS annotations(1415)

**b**

| Gene | VariationType | P | | | |
|------|---------------|---|---|---|---|
| | | | — Univariate | — Multivariate | |
| ERBB2 | nonsilent | 0.2100 | | | |
| | | 0.0690 | | | |
| ERBB4 | nonsilent | 0.3960 | | | |
| | | 0.2150 | | | |
| JAK1 | MS | 0.1700 | | | |
| | | 0.0730 | | | |
| PIK3R1 | nonsilent | 0.0097 | | | |
| | | 0.0021 | | | |
| PIK3R1 | FS/NS | 0.0110 | | | |
| | | 0.0145 | | | |
| RB1 | nonsilent | 0.3000 | | | |
| | | 0.8910 | | | |
| Tumor Grade | clinical | 0.00043 | | | |
| Node Positivity | clinical | 5.24E-13 | | | |
| Tumor Size >5cm | clinical | 6.67E-10 | | | |

Hazard Ratio

0.25    0.50    1.0    1.41    4.0

**Figure 7.**

Figure 8.