# The prognostic effects of somatic mutations in ER-positive breast cancer

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

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

* These authors contributed equally.
† Corresponding authors. matthew.ellis@bcm.edu; elaine.mardis@nationwidechildrens.org

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

**Abstract**

DNA from primary breast cancer samples from 625 postmenopausal (UBC-TAM series) and 328 premenopausal (MA12 trial) hormone receptor-positive (HR+) patients were subjected to targeted sequencing of 83 genes to determine interactions between somatic mutation and prognosis. Independent validation of prognostic interactions was achieved using data from the METABRIC study. Previously established associations between MAP3K1 and PIK3CA with luminal A status/favorable prognosis and TP53 mutations with Luminal B/non-luminal tumors/poor prognosis were observed, validating the methodological approach. As observed in UBC-TAM, *NF1* frame-shift nonsense *(FS/NS)* mutations were also a METABRIC-validated poor outcome driver. For MA12, poor outcome associated with PIK3R1 mutation was also reproducible in METABRIC. DDR1 mutations were strongly associated with poor prognosis in UBC-TAM despite stringent false-discovery correction (q=0.0003). In conclusion, uncommon recurrent somatic mutations should be further explored to create a more complete explanation of the highly variable outcomes that typifies ER+ breast cancer.

**Introduction**

While recent genomic studies have provided a comprehensive catalog of genes that accumulate somatic point mutations and small insertions/deletions (indels) in estrogen receptor-positive (ER+) breast cancer, there remains considerable uncertainty as to how these newly discovered mutations relate to disease outcomes [1, 2, 3]. Most genomic discovery cohorts were neither uniformly treated nor followed long enough. For ER+ disease in particular, prognostic studies require prolonged observation since relapses often occur after 5 years [4]. Uniform treatment was a feature of a whole genome sequencing study of samples accrued from a neoadjuvant aromatase inhibitor (AI) clinical trial for ER+ clinical stage 2 or 3 disease, although only short-term anti-proliferative response to AI were reported. This investigation identified that mutations in *MAP3K1*, a tumor suppressor gene involved in stress kinase activation, were associated with indolent biological features and low proliferation rates [5]. The resulting hypothesis was that *MAP3K1* mutation would be associated with favorable outcomes. In contrast, *TP53* mutations associated with poor prognosis features and high proliferation rates.

To more comprehensively address the relationships between somatic mutations and outcomes in ER+ breast cancer, we developed an approach to detect somatic mutations in DNA isolated from formalin fixed tumor blocks that were over 20 years old. After curating existing mutational data from breast cancer genomics discovery studies (Supplementary Data 1), 83 genes were chosen for analysis (Supplementary Table 1). We applied DNA hybrid capture, sequencing and somatic analysis to three ER+ breast cancer discovery cohorts with contrasting clinical characteristics: An older cohort treated with adjuvant tamoxifen and no chemotherapy (UBC-TAM series[6]), a premenopausal cohort uniformly treated with chemotherapy and randomized to tamoxifen versus observation (NCIC MA12 clinical trial[7]); and a third mixed cohort that was used only to expand the mutational landscape analysis (POLAR) (Supplementary Table 2). An analytical pipeline was developed to identify somatic variants while compensating for the lack of matched normal DNA, which is generally unavailable in the setting of older formalin-fixed tumor material. Somatic mutations were analyzed for association with standard clinical variables, wherein mutated *TP53* and *MAP3K1* served as *a priori* hypotheses for poor and good outcome, respectively. Additional objectives were to identify new mutational hotspots, assess interactions with PAM50-based intrinsic subtypes and to determine mutation frequencies for therapeutic targets. Validation was possible by comparing our results to those in cBioPortal where the genes sequenced in the METABRIC cohort overlapped with the 83 genes investigated in the study described herein.

**Results**

Sequencing and final study cohorts

University of British Columbia Tamoxifen Series (UBC-TAM): These cases were drawn from a well-annotated cohort of patients treated with adjuvant tamoxifen without chemotherapy [6]. A total of 625 of 632 (98.8%) patient samples that fully met study criteria passed a minimum sequencing quality cutoff of at least 80% of targeted exonic bases covered at greater than 20X (mean coverage: 133X) with other quality metrics described in the

111 supplementary data (Supplementary Figure 1-5 and Supplementary Data 2). Mean depth was correlated with
112 input DNA and negatively correlated with time since diagnosis (approximate age of sample) and duplication
113 rates were negatively correlated with input DNA and positively correlated with sample age. However, despite
114 these trends, overall metrics were excellent with an average of 135.8X coverage and 3.0% duplicate rate
115 despite the generally low input amounts and old sample age. The final patient population had an average age
116 of 67 at diagnosis (range: 40-89+).  All were treated with five years of adjuvant tamoxifen, and were primarily
117 postmenopausal, grade 2 or 3 cancers, of ductal histologic subtype (Supplementary Table 2). All were ER+
118 (>1% cells positive by IHC) and at least 88.6% were clinically HER2- (13/625 unknown).  A subset of 463 of
119 these patients had PAM50 subtyping data available from a previous study [6]. The median follow up in the cohort
120 examined was 25 years and one month.
121
122 NCIC-MA12 Trial cohort: These cases were drawn from a clinical trial in premenopausal women treated with a
123 standard adjuvant chemotherapy regimen and randomized to tamoxifen versus observation. A total of 459
124 patient samples passed the minimum sequencing quality threshold (mean coverage: 200X), of which 328 were
125 hormone receptor positive (HR+; >1% cells positive for ER or PR by IHC), and only the HR+ cohort were
126 included here for most analyses. The majority were premenopausal (mean age of 45). All patients received
127 chemotherapy, and 48% were treated with 5 years of adjuvant tamoxifen. A subset of 255 of these patients
128 had PAM50 subtyping data available. The median follow up in the cohort examined was 9.7 years
129
130 POLAR cohort: This patient series was a case-control study of ER+ (>1% cells positive by IHC) breast tumors,
131 175 of 194 (90.2%) patient samples passed minimum sequencing quality thresholds (mean coverage: 75X). A
132 case was defined as any patient who relapsed during follow-up, and controls were defined as lacking relapse
133 through a similar follow-up duration. Based on these definitions, there were 91 cases and 84 controls. Of the
134 cases, 43 were early relapses (<5 years since diagnosis) and 48 were late relapses (>5 years). Patients were
135 only included if they received adjuvant endocrine therapy, but chemotherapy was not an exclusion criterion,
136 nor was menopausal status. Because the POLAR study was a case-control design, outcome data could not be
137 easily integrated into prognostic analysis. Therefore, these cases were used in the mutation landscape and
138 hotspot analyses only.
139
140 Across the three cohorts, there were 1,259 patient samples that passed minimum sequencing quality
141 thresholds and 1,128 of these were ER+ (UBC-TAM and POLAR) or ER and/or PgR+ (HR+) (MA12).
142
143 Variant calling and filtering
144
145 A total of over 62 million variants were identified in UBC-TAM. After extensive filtering against a set of nearly
146 70,000 unmatched normal samples and manual review to eliminate common polymorphisms and false
147 positives (see methods), 1,991 putative somatic variants were identified (0 to 26 variants per patient). A set of
148 1,693 mutations was defined as the "non-silent" set for further analysis that excluded sequencing variants in
149 splice regions (except proximal splice site), RNA genes (except *MALAT1*), UTRs, introns, and all silent
150 mutations. Finally, a set of 408 frameshift or nonsense mutations was defined. The same filtering method was
151 applied to both the POLAR and MA12 datasets. A total of 540 putative somatic mutations (436 non-silent, 145
152 FS/NS) were identified in POLAR, and 2,104 (1,753 non-silent, 610 FS/NS) in MA12. Full details on these
153 variants are included in Supplementary Data 3 and summarized for key genes in Supplementary Figure 6.
154
155 Mutation landscape analysis
156 In 1128 samples passing quality control standards, considering only non-silent mutations, 17 genes were
157 mutated at a rate greater than 5%, and 6 at a rate greater than 10%; *PIK3CA* was the only gene mutated in
158 greater than 20% of samples (**Figure 1A**). The order from most recurrent to least for the 10 most frequently
159 mutated genes was: *PIK3CA* (41.1%), *TP53* (15.5%), *MLL3* (13.4%), *MAP3K1* (12.0%), CDH1 (10.5%),
160 MALAT1 (10.0%), GATA3 (9.1%), MLL2 (8.7%), ARID1A (7.2%), and BRCA2 (6.6%). This list correlates well
161 with previously reported recurrently mutated genes. For example, the top 4 most significantly mutated (non-
162 silent) genes in the ER+ subset of TCGA breast project[3] were *PIK3CA* (24.0%), *TP53* (14.6%), *GATA3* (8.6%)
163 and *MAP3K1* (6.1%). Considering METABRIC ER+ patients, the most recurrently mutated genes were PIK3CA
164 (~46%), TP53 (~21%), GATA3, MLL3, CDH1, and MAP3K1 (all ~12-14%) demonstrating slightly higher but

165 very similar frequencies. The overall average non-silent mutation frequency was estimated as 1.6 per MB of
166 coding sequence (range: 0.5 to 5.8 mutations per MB, excluding samples with no mutations called). In order to
167 determine whether mutations in any gene pair were mutually exclusive or co-occurring in this dataset, a
168 pairwise Chi-squared or Fisher's exact test was performed. Mutations in PIK3CA and MAP3K1 were
169 significantly more likely to co-occur (after BH FDR correction) in TAM dataset, and were near significance in
170 MA12 although not after correction (p = 0.08). These results are summarized in Supplementary Data 4.
171
172 Hotspot analysis
173
174 As anticipated [8], mutations in *PIK3CA* at *E542K*, *E545K*, and *H1047R* were highly recurrent in this study with
175 69/1259 (5.5%) E542K, 104 (8.3%) E545K, and 181 (14.4%) H1047R mutations (Supplementary Figure 6C).
176 Mutations in the ligand binding domain of *ESR1* (1.1%) were extremely rare [3, 9, 10] (Supplementary Figure 6A).
177 To uncover novel hotspots in these data, both Chi-squared and Fisher's exact tests were performed using
178 mutation frequencies from previous sequencing studies as the expected values (see Methods for definition of
179 multi-study MAF file) (Supplementary Table 3). The most notable novel finding was in *CBFB* (**Figure 1B**). At
180 least 6 different genomic alterations were observed in 15 patients (Supplementary Data 3) that affected the
181 donor splice site of exon 2. Manual review of this splice site identified at least two additional patients with
182 evidence for mutations at this location. The predicted effect of these mutations is skipping of exon 2 or
183 alternate donor site usage, each likely resulting in loss-of-function of the *CBFB* protein. Additional splice site
184 mutations were observed at the exon 2, exon 4 and exon 5 acceptor sites of *CBFB*.  ErbB2 exhibited the
185 anticipated profile of activating mutations from earlier publications[11] with 22/1259 (1.7%) samples harboring
186 known activating mutations and another 6 variants of unknown significance in the kinase domain or at the
187 S310 residue (**Figure 8C**).
188
189 Somatic mutation association with PAM50-based intrinsic subtype
190
191 PAM50 intrinsic subtype calls were obtained from previously published analyses to compare to their mutational
192 profiles for UBC-TAM and MA12 (HR+ only) studies.  In both studies about half the patients had luminal A
193 tumors.  However, the MA12 cohort had a higher proportion of non-luminal subtypes, with 19.8% HER2-E and
194 6.6% basal and fewer luminal B tumors (25.1% versus 42.4%) (**Figure 2A-B**). As expected, patients with the
195 HER2-E intrinsic subtype were enriched for HER2+ve status compared to other subtypes (Fisher's exact test
196 p<0.0001). Of interest, in the HER2-enriched group there were 51 tumors that were not HER2 amplified and of
197 these 4 were HER2 mutant (~8%), indicating that HER2 mutation could be an occasional explanation for a
198 HER2-E subtype assignment in the absence of HER2 amplification.  For NF1 FS/NS mutations, there was also
199 a statistically significant association with the HER2-E subtype (P=0.002) (supplementary Figure 7B, also
200 supplementary data 5).  Notably NF1 non-silent mutations were enriched in the HER2-E non-HER2 amplified
201 subgroup, where they were present in 8/51 cases (16%).  Compared to the frequency in all other subtypes
202 12/582 (2%), this enrichment was significant (Fishers exact test p<0.0001) (Supplementary 7A right panel).
203 This association could be reproduced in the METABRIC data with an NF1 non-silent mutation incidence in the
204 HER2-E non HER2 amplified group of 8/80 (10%) versus 35/1283 (3%) in the rest of the subtypes (p=0.003)
205 (Supplementary 7A right panel).   Age density plots by subtype serve to emphasize the large difference in the
206 median age between the two sample cohorts (43 versus 65), and also the influence of age with respect to the
207 intrinsic subtype incidence. Namely, in the younger MA12 cohort, there is a younger peak incidence with basal-
208 like breast cancer than Luminal A disease (**Figure 2D**). In contrast in the older UBC-TAM cohort, an influence
209 of age on intrinsic subtype was not observed (**Figure 2C**).   Relationships between intrinsic subtype and
210 mutation patterns were also explored, classifying mutation positive status as "non-silent", "missense",
211 nonsense/frame-shift (FS/NS) or FS/NS+splice site (Supplementary Data 5).  The FDR corrected p-value (q-
212 value) took into account that 83 genes were examined.  However, this level of false discovery detection could
213 be viewed as overly conservative in an exploratory analysis.  Therefore, any gene mutation with q-value
214 association of <0.2 was therefore considered reportable for the purposes of subsequent validation efforts [12, 13,
215 14].  For MA12, non-silent TP53 mutation was highly subtype-associated because of the very high incidence in
216 non-luminal versus luminal subtypes.  PIK3CA and MAP3K1 mutations were associated with Luminal A
217 disease in both cohorts (Supplementary Figure 7B). Finally, there was a strong association between Luminal B
218 status and non-silent (Supplementary Figure 7A) as well as FS/NS mutations in GATA3 (Supplementary Data

3

219   5, q value = 0.006) for MA12 (but not UBC-TAM).  GATA3 mutations were present in 28-30% of Luminal B
220   cases and less so in luminal A cases (5%). Considering q values of <0.2 the associations between FS/NS and
221   non-silent mutations in ATM and Luminal B tumors in MA12 (8-13%) suggests that ATM disruption is also a
222   possible luminal B driver (Supplementary Figure 7C), at least in younger women (MA12). Relationships
223   between age and mutation incidence were therefore also explored (Supplementary Figure 7D), with the finding
224   that both ATM mutation and GATA3 mutations were associated with an earlier age of onset within the luminal
225   B category (**Figure 2E and 2F**).  Some ATM mutations are likely to be germline (see discussion below), which
226   could partially explain the association with younger age.

227
228   <u>Survival analysis according to somatic mutation.</u>
229
230   For the UBC-TAM Series (**Figure 3A**) univariate analysis, favorable prognostic associations for breast-cancer-
231   specific survival (BCSS) were detected for non-silent mutations in *MAP3K1*, *ERBB3*, XBP1 and PIK3CA
232   (**Figure 3B**, Supplementary Data 6). Adverse prognostic effects were observed for non-silent mutations in
233   *DDR1* and *TP53*, as well as for frame-shift and nonsense (FS/NS) mutations in NF1. An analysis for
234   recurrence free survival (RFS) produced similar results, except for ARID1B, which was marginally associated
235   with more favorable outcome. A multivariate Cox model was applied to put each gene in the context of clinical
236   parameters (grade, tumor size and node status).  These analyses indicated that the prognostic effects of non-
237   silent DDR1, PIK3CA, GATA3 FS/NS, TP53 and MAP3K1 mutations were independent of grade and
238   pathological stage (**Figure 3C**).  Multiple correction testing, yielded DDR1 as the only gene that remained
239   significant with a q-value of 0.0003 (Supplementary Data 5).  For the MA12 clinical trial cohort (**Figure 4A**) we
240   focused on overall survival associations, as this was the primary endpoint of the study and the most robust
241   endpoint.  A number of rarely mutated genes were associated with poor outcome in univariate analysis as
242   displayed in **Figure 4B**. Multiple testing corrections indicated none of these findings could be considered
243   significant [12, 13, 14].  However, in multivariate analysis, based on the uncorrected p value, the prognostic effects
244   of mutations in ErbB2, ErbB4, LTK FS/NS, MAP3K4, PIK3R1, RB1, RELN and TGFB2 were independent of
245   pathological stage and grade (**Figure 4B**).

246
247   <u>Verification of Prognostic effects of Mutations in METABRIC data.</u>
248
249   While few genes were significant in univariate analysis after multiple testing correction, their identification
250   provides valuable hypotheses for further testing and validation. We therefore sought additional data in the
251   public domain to further assess the uncorrected p value-based findings in our data set.  The METABRIC
252   consortium have reported somatic mutations in cBioPortal [15] with co-reported detailed hormone receptor
253   status, age at diagnosis (median age=64 years for ER+ patients), mean follow up of >8 years and disease-
254   specific (breast-cancer-specific) outcome [16, 17]. This data set provided the opportunity to conduct a validation
255   exercise for overlapping genes in the two data sets.  For the UBC-TAM series (**Figure 3**), 9 genes with a
256   univariate p value of <0.05 were brought forward for validation (**Figure 5**).  Of the 6 overlapping genes also
257   examined in METABRIC, consistent prognostic effects independent of clinical variables were observed for non-
258   silent mutations in three genes, *MAP3K1* (favorable), *TP53* (unfavorable) and *NF1* FS/NS mutations
259   (unfavorable).  In order to maintain coherence in discovery and validation patient cohorts, a similar analysis
260   was carried out restricting the patient pool to postmenopausal patients only. No significant variation in hazard
261   ratio for candidate genes where observed (Supplementary Table 4). For the MA12 series (**Figure 4**), 5 shared
262   genes were identified with univariate p values of <0.05, yet only *PIK3R1* mutations (non-silent or FS/NS)
263   showed consistent adverse prognostic effects (**Figure 6**).  The Kaplan Meier survival plots for the consistent
264   adverse prognostic effects of *NF1* FS/NS (TAM vs METABRIC) and non-silent *PIK3R1* (MA12 vs METABRIC)
265   mutations are illustrated in **Figure 7A-D**. Copy number aberrations and chromosomal instability have been
266   associated with prognosis across multiple cancer types, including ER-positive (ER+) breast cancer[16, 18, 19]. To
267   gauge the confounding nature of commonly amplified genes in breast cancer, we further performed
268   multivariate analysis on the candidate genes with cases of amplification of MYC, FGFR1, CCND1 and ERBB2
269   (Supplementary Table 5). We did not observe a significant change in the hazard ratio reported in Figure 5B
270   and 6B).

271
272   <u>Prognostic interactions between PIK3CA and MAP3K1.</u>

273
274 Since PIK3CA and MAP3K1 mutations co-associate, the combined effect of non-silent mutations in these
275 genes was examined. Patients with tumors exhibiting both genes mutated have a more favorable clinical
276 course than either singly mutant cases or cases without either gene mutated. While the prognostic effects
277 were strongest in the UBC-TAM series, this result was also reproduced in the METABRIC data (**Figure 7E-F**).
278
279 <u>Mutation Analyses for Uncommon Targetable Kinases.</u>
280
281 Of the 83 genes analyzed, at least 8 are directly targetable with small molecules or antibodies that are either
282 FDA approved or in late-stage development (**Figure 8**). Pre-existing data on these mutations is summarized
283 (Supplementary Data 7). PIK3CA is not further discussed here, since the mutation spectrum is well described
284 and large therapeutic studies are already underway. An examination of the 23 mutations in ERbB2 revealed
285 locations that were, as expected, clustered in 2 major domains, with 2 of 23 having extracellular domain
286 mutations at residue 310 and 21 of 23 having kinase domain mutations between residues 755-842 [11, 20]. To
287 further investigate the preliminary finding of an adverse prognostic effect for ErbB2 mutation in the MA12
288 series, an examination of the METABRIC data indicated that known activating mutations in ErbB2 were
289 associated with a near significant adverse effect (HR=1.71, P=0.075) (Supplementary Figure 8).
290
291 For ERBB3, 2 known-activating mutations were identified (V104L and E928A)[21]. The DDR1 kinase domain
292 mutation, R776W, is possibly homologous to EGFR hot spot mutation L858R, but the remaining DDR1 variants
293 are of unknown significance. For the mutations in JAK1, 3 of 12 are loss of function mutations (frame shift or
294 non-sense) and the S816* mutation has been reported in a lung adenocarcinoma sequencing data set [22]. The
295 loss of function mutations in JAK1 have been shown to associate with immunotherapy resistance [23, 24]. A few
296 mutations identified in ERBB4, MET, and PDGFRA have been previously reported but those reported here
297 have not been functionally tested.
298
299 **Discussion**
300
301 The strength of this investigation includes the prolonged follow up, controlled adjuvant treatment and the
302 relatively large number of genes and patients studied. Weaknesses include the lack of treatment prediction
303 because endocrine treatment in UBC Tam was uniform but not randomized. In MA12 the use of tamoxifen was
304 randomized, but the numbers were too small to examine treatment interactions. The landscape of recurrently
305 mutated genes in ER+ breast cancer observed in this study is consistent with reports where matched germline
306 samples were available, indicating that our variant filters were effective for somatic mutation detection in a
307 research setting. Overall, mutation frequencies were higher in our cohort (e.g., for *PIK3CA*, *MLL3*, *MAP3K1*)
308 than the TCGA cohort, but were also lower for a few specific genes (e.g., *TP53* and *GATA3*). Due to higher
309 sequencing data coverage of recurrently mutated target genes than TCGA and the use of a different hybrid
310 capture reagent, we were likely able to detect mutations that were missed with lower-depth exome or whole
311 genome sequencing data. Differences in patient populations may also be a factor. Frequencies were much
312 closer to reported values for METABRIC which also used a targeted sequencing approach. It is also possible
313 that in some instances we overestimated somatic mutation frequency, due to the lack of matched normal
314 samples and imperfections in our germline polymorphism filtering. In particular, a significant number of *BRCA1*,
315 *BRCA2*, and *ATM* mutations are likely *de novo* germline mutations that we would not be able to easily
316 distinguish from somatic mutations. Of the 117 non-silent *BRCA1/2* mutations observed (from 110/1128
317 patients across all 3 cohorts; 7 patients had two hits) 74 were observed at a VAF greater than 40% and 31
318 were greater 60%. Additionally, of the 61 non-silent *ATM* mutations (from 58/1128 samples; 3 samples had 2
319 hits) 39 had VAF greater than 40% and 18 had VAF greater than 60 (Supplementary Data 9). Variants with
320 VAFs this high are less likely to be somatic given the general expectation of impure tumor samples and
321 heterozygous mutations. Indeed, the VAFs for *BRCA1/2* and *ATM* non-silent mutations (mean=46.0%) were
322 significantly higher than for other genes (mean=36.7%, p=5.92e-09). Even when considered separately, the
323 VAFs for *BRCA1* (mean=46.6%), *BRCA2* (mean=43.8%) and *ATM* (mean=48.2%) were significantly higher
324 than the other genes (p=0.002, p=0.0015, and 5.27e-5 respectively). Among the *BRCA1/2* variants, there were
325 8 known pathogenic (ENIGMA expert reviewed) mutations according to a search of the BRCA Exchange
326 database (http://brcaexchange.org, Nov 12, 2017) and another 37 assumed pathogenic (FS/NS) mutations. Of

5

the remaining, 4 were benign according to expert review (ENIGMA), and 8 benign, 15 likely benign and 45 variants of unknown significance according to all public sources. Out of the 61 *ATM* variants queried in ClinVar, 4 were designated as pathogenic, 3 were pathogenic/likely pathogenic, and 2 were likely pathogenic. Another 7 were frameshift mutations and assumed pathogenic. Additionally, 23 variants had uncertain significance, 8 variants had conflicting interpretations of pathogenicity (any combination of benign, likely benign, or uncertain significance), and the remaining 14 variants had no data. *ATM* variants were also queried in the Leiden Open Variation Database (LOVD) [25], which identified 1 variant that affects function (designated as likely pathogenic by ClinVar), 10 variants with unknown effect, and 1 variant that probably does not affect function (uncertain significance in ClinVar). The remaining variants had no data in LOVD. Given these complexities the prognostic effects of somatic versus germline BRCA1/2 and ATM mutations remain unresolved, however attention should clearly be paid to therapeutic strategies for these patients. The ATM findings deserve a particular highlight because of the younger age/luminal B association and the current lack of studies devoted to this population.

The discovery of a novel recurrent *CBFB* (core binding factor subunit beta) splice site mutation in this cohort illustrates a limitation of exome capture reagents. The affected bases in exon 2 of CBFB display reduced sequence coverage, possibly due to high GC content, in the breast TCGA exome dataset (Supplementary Figures 9-10). This site was mutated in at least 1.5% of ER+ breast cancers sequenced, bringing the overall rate of CBFB mutations to nearly 6%, which should drive further investigation of this gene in ER+ breast cancer pathogenesis. *CBFB* functions as a subunit in a heterodimeric core binding transcription factor that interacts with *RUNX1*[26]. Consistent with this model, *CBFB* mutants were mutually exclusive from *RUNX1* mutants in this cohort with only a single sample harboring non-silent mutations in both *CBFB* and *RUNX1*.

The UBC-TAM and MA12 studies revealed different lists of potentially prognostic mutations. Prognostic effects are likely to be strongly affected by the use of systemic therapy as well as by patient age at diagnosis. The UBC-TAM series is the simplest study to interpret from a drug resistance perspective since the only systemic therapy was tamoxifen. Thus, the consistent adverse effect of NF1 FS/NS mutation on prognosis is intriguing as this result is consistent with results from an *in vitro* screen for tamoxifen resistance[27]. Understanding why only FS/NS mutations predict poor outcome, rather than missense or other non-silent mutations, will require further investigation. The association with the HER2-E, non-HER2 amplified subset with non-synonymous NF1 mutations was observed in both the discovery and validation (METABRIC) data sets. It is a logical proposition that mutations that activate RAS, like NF1 mutation, could create a tumor with a similar transcriptional phenotype as some HER2 amplified breast cancers. PIK3R1 mutation also emerged as a consistent poor prognosis mutation from the MA12 analysis, with validation in METABRIC. The proposed favorable prognostic effects of PIK3CA mutation were observed in the UBC-TAM series, but were not found to be independent of stage and grade, and PTEN mutations were neutral.

According to our validation results, NF1, PIK3R1, PIK3CA and TP53 are therefore likely to be prognostic drivers that are independent of clinical variables. In postmenopausal women treated with adjuvant endocrine therapy, DDR1, PRKDC and XBP1 should be further studied and of these DDR1 is the strongest candidate because it was significant despite strict false discovery correction. DDR1 is a collagen-binding receptor expressed in epithelial cells that stabilizes E-cadherin–mediated intracellular adhesion[28]. *DDR1* mutations also occur in endometrial cancer[29], acute leukemia[30] and lung cancer[31]. Loss of DDR1 (DDR1-null mice) produces hyper-proliferation and abnormal branching of mammary ducts, suggesting DDR1 is a breast tumor suppressor[32]. Mutations in PRKDC will potentially produce a defective ATM response/low ATM levels [33] which is interesting in the context of the finding herein that ATM mutations are a potential luminal B driver gene. The significance of a defective ATM pathway as a cause of endocrine resistance is highlighted by the recent finding that dysregulation of the MutL complex (MLH1, PMS1 and PMS2) causes failure of ATM/CHK2-based negative regulation of CDK4/6 [34]. Prognostic candidate mutations revealed by the MA12 analysis were different from the UBC TAM series, likely reflecting the different patient profiles and adjuvant treatments illustrated in **Figure 2**. The prognostic effects of mutations ERBB2, ERBB4, JAK1, LTK, MAP3K4, MET, PDGFRA, RB1, RELN, TGFB2, all await further study with even larger sample sizes.

A limitation of this study is that the mutation datasets we generated for UBC-TAM and MA12 cohorts lack comprehensive assessment of copy number signatures that have been associated with prognosis in ER+

6

breast cancer[16, 18, 19]. While multivariate analysis considering key CNVs did not appear to affect our prognostic associations, future studies may be needed to completely understand the interplay between simple and large-scale variation for prognostic prediction. Another limitation to this study was the heterogeneity in the datasets in terms of age, treatment, and other factors that limited direct comparison and made validation with METABRIC somewhat challenging. The collection of sufficiently large, uniformly treated populations with long-term follow-up for discovery and validation remains a challenge that must be addressed to fully characterize the prognostic significance of somatic mutations, especially low-frequency mutations.

In conclusion, we have successfully utilized clinically well-annotated, uniformly treated patient samples using DNA from archival material greater than 20 years old without a matched normal to explore the prognostic effects encoded by the mutational landscape of ER+ breast cancer. We were able to confirm our prospective hypothesis from our earlier studies [5] that MAP3K1 is associated with indolent disease and TP53 with adverse outcomes. We also associated NF1 FS/NS mutations with strong adverse effects on prognosis. Similarly, PIK3R1 mutations were associated with an adverse prognosis, in contrast to PIK3CA mutation which were weakly favorable. This suggests somatic mutations in these two physically interacting gene products are not biologically equivalent with respect to PI3 kinase pathway activation and resistance effects. The possibility that the long tail of low frequency mutation events in luminal type breast cancer may harbor multiple molecular explanations for poor outcomes should spur new collaborative efforts to thoroughly screen thousands of properly annotated cases. Only after these iterative efforts of proposing and confirming candidates will a clinically useful and comprehensive somatic mutation-based classification of ER+ breast cancer emerge. In the meantime, functional studies should be pursued to understand the biological effects of low frequency somatic mutations, prioritizing these studies according to whether the mutations are driving an adverse prognostic effect and whether their disruption creates a therapeutic vulnerability.


**Methods**

For the UBC-TAM series, an institutional review board approved study was based on formalin-fixed paraffin embedded (FFPE) primary tumor blocks from 947 female patients diagnosed with estrogen receptor positive invasive breast cancer in the province of British Columbia in Canada between 1986 and 1992[6, 35, 36, 37]. The sample flow and analysis are provided in a REMARK summary (**Figure 3A**). DNA was isolated from tumor-rich regions using the Qiagen blood and tissue kit, which yielded sufficient DNA in 645 samples, of which 625 met all study criteria and had sufficient sequence coverage. Similarly, approved studies provided 194 and 454 HR+ patient samples for the POLAR and MA12 (**Figure 4A**) cohorts. A total of 175 POLAR and 459 (328 HR+) MA12 samples yielded sufficient DNA and had sufficient sequence coverage for analysis. Detailed descriptions of the patient data sets are provided in Supplementary Table 3. A meta-analysis of six existing published large-scale breast cancer sequencing studies [1, 2, 3, 5, 38, 39] was performed to identify genes with recurrent coding region somatic mutations in breast cancer (Supplementary Data 1). Additional drug targets[40] and genes with relevance to breast cancer from targeted sequencing[41], copy-number studies[16] or knowledge relating to somatic or germline mutations (e.g., *BRCA1*, *BRCA2*, *ERBB2*, *ESR1* and *PRLR*) were also included. This resulted in a final list of 83 breast-cancer-related genes (Supplementary Table 1). These genes were targeted comprehensively with 3,029 complementary probes for hybridization-based enrichment (Supplementary Data 8). Sequencing libraries were constructed, hybridized with capture probes, multiplexed and run on a single flow cell with up to 96 samples per pool per lane yielding approximately 375 Mb of DNA sequence per sample from an Illumina HiSeq paired end 2 X 100bp (TAM) or 2 X 125bp (POLAR, MA12) sequencing run following manufacturer's protocols.

Variant calling was performed with the Genome Modeling System as previously described[42]. Specifically, sequence data were aligned to reference sequence build GRCh37 using BWA[43] and de-duplicated with Picard. SNVs and indels were detected using the union of samtools[44] and VarScan2[5] and annotated using Ensembl version 70. Variants were restricted to the coding regions of targeted genes and filtered for false positives and germline polymorphisms against a database of nearly 70,000 unmatched normals from the ExAC consortium[45], 1000 Genomes[46], NHLBI exomes[47] and TCGA data sets[3, 48]. A binomial probability model was then applied to the variants using VAF and total coverage to determine a log-likelihood ratio of being a somatic variant as

previously described[49] (See Supplementary Methods). After filtering, all remaining variants were manually reviewed. To ensure that variants of known clinical relevance were not missed by automated variant calling approaches, a knowledge-based variant calling strategy was performed focused on the mutations in the Database of Curated Mutations[50].

Patient groups were defined by mutation status or truncating mutation status for each gene. Fisher's exact and Chi-squared tests were used for hotspot analysis, mutual exclusivity or co-occurrence, and other categorical clinical statistics (e.g., mutation status vs. intrinsic subtype) as appropriate. Univariate Kaplan-Meier and Cox survival analyses were performed for breast-cancer-specific survival (BCSS), relapse free survival (RFS), or overall survival (OS) with non-silent or truncating mutation status as a factor. Significant survival differences between the groups were determined by log rank (Mantel-Cox) test. The Benjamini-Hochberg method was performed for multiple testing corrections to report the false discovery rate adjusted p-value (q-value). A multivariate Cox proportional hazard model was fitted to BCSS and RFS separately on gene mutation status, node status, grade and tumor size and adjusted hazard ratios were calculated with Wald test p-values. All statistical analyses were performed in the R statistical programming language with core, 'survival' and 'multtest' libraries. Genomic visualizations were created with ProteinPaint[51] and GenVisR[52].

## Data Availability

All mutation calls are made available as a MAF file with this publication. The raw sequence data from UBC-TAM patients are available in the database of Genotypes and Phenotypes (dbGaP) under accession number [dbGAP:phsxxxx.x]. Raw sequence data from MA12 and POLAR could not be deposited in public repository due to patient consent issues and complexities of institutional certification. However, these data are available from the authors (contact Obi Griffith and Matthew Ellis). Primary clinical outcome data for UBC-TAM and MA12 can be made available to qualified researchers through application to the Canadian Cancer Trials Group. Primary clinical outcome data for POLAR can be made available to qualified researchers through application to Mitch Dowsett at the Ralph Lauren Centre for Breast Cancer Research.

## Contributions

O.L.G., T.O.N., M.J.E., and E.R.M. designed the experiments; N.C.S, M.A., M.G., J. K., C.A.M., K.K., 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., and R.K.W. generated the sequencing data. T.O.N., B.Y., M.D., S.L., and D.V. orchestrated the sample pipeline, M.A., O.L.G. and N.C.S. prepared the figures and tables. O.L.G., N.C.S., M.A., J.L., and D.T. provided statistical analysis. S.M.K., R.B., and E.C.C. provided functional annotations. T.O.N. provided pathology analysis. M.J.E., N.C.S., M.A., and O.L.G. wrote the manuscript. E.R.M., T.O.N., and M.D., critically read and commented on the manuscript.

## Conflict of Interest

**Figure Legends**

**Figure 1. Mutation recurrence and novel splice site mutation**
A) The overall mutation recurrence rate ranged from 41.1% of samples for *PIK3CA* to 0.0% for *PIN1*. The figure depicts non-silent mutations for all 1128 patients for the top 17 most recurrently mutated genes (>5% recurrence). If a patient had multiple mutations it is colored according to the "most damaging" mutation following the order presented in the Mutation Type legend (vertical color bar). Mutations per MB were calculated using the total number of mutations observed over the total exome space corresponding to the tiled space from "SeqCap EZ Human Exome Library v2.0". A correction factor was applied to account for genes not assayed using the expected number of additional mutations based on ER+ TCGA data. The coverage histogram (top sidebar) shows the percent of targeted exonic bases with at least 20X, 30X and 40X coverage. B) Mutation recurrence frequencies (amino acid level) in this study were compared to previously reported mutation frequency from a multi-study MAF file of six reported breast cancer sequencing studies (Supplementary Data 1). An entirely novel mutation "hot spot" was discovered affecting the exon 2 splice (donor) site of *CBFB* in at least 15 patients. Six different single nucleotide substitutions, insertions and deletions were observed, all affecting either the first or second base of the donor splice site. These mutations were most likely missed in previous studies because of a lack of sequencing coverage due to the GC-rich nature of exons 1 and 2 of *CBFB* (Supplementary Figures 9-10). Such mutations are predicted to significantly alter the canonical donor site and result in either alternate donor usage or skipping of one or more exons of *CBFB*.

**Figure 2. Cross-cohort age and subtype analysis**
A-B) Percentage composition of samples by intrinsic subtype of the tumor in the two discovery cohorts for UBC-TAM (A) and MA12 (B) cohorts. C-D) Age-density plots for patients categorized by intrinsic subtype in UBC-TAM (C) and MA12 (D) cohorts. The overall median age shows that UBC-TAM is constituted mostly of post-menopausal patients (median age=65), in contrast to MA12, which has younger patients (median age=43). E-F) Younger luminal B subtype patients harbor GATA3 (E) and ATM (F) mutations in the combined set of UBC-TAM and MA12 Luminal B cases (median age=52, p=0.01; median age=58, p=0.03 for GATA3 and ATM respectively).

**Figure 3. Candidate discovery from UBC-TAM cohort and prognosis evaluation**
(A) DNA was extracted from tumor specimens from 947 patients with ER+ breast cancer treated with tamoxifen monotherapy for 5 years. 632 samples with adequate yield were sequenced for 83 genes known to be recurrently mutated or breast cancer relevant. A total of 625 samples passed minimum quality checks and were sequenced to an average of 135.8X coverage. A total of ~62 million variants from the reference genome were identified. Extensive filtering and manual review reduced this list to 1,991 putatively somatic variants. Survival analysis was applied to non-silent and truncating gene mutation status versus disease outcome (relapse or breast-cancer-specific death). In addition, mutations were analyzed for novel hotspots, patterns of mutual exclusivity or co-occurrence and association with clinical variables. (B) Forest plot of impact of mutations in candidate genes, identified using the UBC-TAM population, on breast-cancer-specific survival (red) and recurrence-free survival (blue). The variant types are characterized based on non-silent or nonsense/frameshift (FS/NS) mutations. The box size is relative to frequency of mutations in the analysis, with larger boxes representing higher incidence mutations. (C) Multivariate forest plot of effect of mutations in UBC-TAM candidate genes on breast-cancer-specific survival when assessed together with clinical factors including Tumor Grade, Node positivity and Tumor Size (>5cm).

**Figure 4. Candidate discovery from MA12 cohort and prognosis evaluation**
(A) DNA was extracted from tumor specimens and 470 samples with adequate yield were sequenced for 83 genes known to be recurrently mutated or breast cancer relevant. A total of 459 (328 HR+) samples passed minimum quality checks and were sequenced to an average of 272.6X coverage. A total of 406 million variants from the reference genome were identified. Extensive filtering and manual review reduced this list to 2104 putatively somatic variants. Survival analysis was applied to non-silent and truncating gene mutation status versus overall survival. (B) Forest plot showing effect of mutation in candidate genes on overall survival (univariate - blue, multivariate - orange), along with the clinical factors used in the multivariate analysis (black),

547  tumor grade, node positivity and tumor size (>5cm). The box size is relative to frequency of mutations in the
548  analysis, with larger boxes representing higher incidence mutations. Note: a few boxes are not shown if their
549  hazard ratios were greater than 4.0.
550
551  **Figure 5. Validation of UBC-TAM candidates in ER+ METABRIC**
552  A) Six out of nine candidate genes from UBC-TAM analysis had mutations reported in the METABRIC cohort.
553  1,060 ER+ samples with breast-cancer-specific survival information were used to test the effect of mutations in
554  the candidate genes on prognosis. B) Forest plot shows effect of mutated candidate genes on breast-cancer-
555  specific survival in METABRIC ER+ cohort with univariate cox proportional-hazard ratio in blue and multivariate
556  in orange.  The clinical factors used in the multivariate analysis, namely tumor grade, node positivity and tumor
557  size (>5cm), are shown in black. The box size is relative to frequency of mutations in the analysis, with larger
558  boxes representing higher incidence mutations. The # cases/CNV column shows the total number of cases
559  with the SNV/Indel variant surrounded by a ring chart indicating the proportion of total cases with CNV
560  alterations.
561
562  **Figure 6. Validation of MA12 candidates in ER+ METABRIC**
563  A) Five out of eleven candidates from MA12 analysis had mutations reported in the METABRIC cohort. 1,415
564  ER+ samples with overall survival information were used to test the effect of mutations in the candidate genes
565  on prognosis. B) Forest plot shows effect of mutated candidate genes, shortlisted based on MA12 mutation
566  analysis, on overall survival in METABRIC ER+ breast cancer patients. Univariate (blue) and multivariate
567  (orange) cox proportional-hazard ratio depict the independent prediction of survival outcomes for the six
568  candidate genes. The box size is relative to frequency of mutations in the analysis, with larger boxes
569  representing higher incidence mutations. The # cases/CNV column shows the total number of cases with the
570  SNV/Indel variant surrounded by a ring chart indicating the proportion of total cases with CNV alterations.
571
572  **Figure 7. Kaplan-Meier plots**
573  A-B) Kaplan-Meier graphs showing the prognostic role of NF1 mutations, separated by variant type – Missense
574  (MUT MS, green), Frameshift/Nonsense (MUT FS/NS, blue) in ER+ breast cancer patients from A) UBC-TAM
575  and B) METABRIC cohort establishing the association between FS/NS mutations in NF1 with poor prognosis.
576  C-D) Kaplan-Meier graph showing the prognostic role of PIK3R1 in C) MA12 and D) METABRIC ER+ breast
577  cancer patients, categorized based on tumors with wildtype (WT, black) or mutated PIK3R1 non-silent
578  mutations (MUT, red). E-F) Kaplan-Meier graph demonstrating co-occurrence of non-silent mutations in
579  MAP3K1 and PIK3CA (red) in E) UBC-TAM and F) METABRIC associates with better survival when compared
580  against tumors with mutations exclusively in MAP3K1 (blue) or PIK3CA (green) or wildtype for both MAP3K1
581  and PIK3CA (black). p, log rank (Mantel-Cox) test p-value.
582
583  **Figure 8. Mutation profiles for selected genes**
584  Mutation frequency plots illustrate all non-silent mutations (TAM, POLAR, and MA12; n=1259) for
585  representative transcripts for several kinase genes of interest. The domains belonging to A) DDR1 (RefSeq ID:
586  NM_013994) and B) JAK1 (NM_002227) are indicated below the schematic diagram of each gene. The ECD
587  (extracellular domain), TM (transmembrane domain), and kinase domain are depicted as green, red, and
588  orange bars respectively for C) ERBB2 (NM_004448), D) ERBB3 (NM_001982), E) ERBB4 (NM_005235), F)
589  MET (NM_000245), and G) PDGFRA (NM_006206). The variant counts across the three datasets for each
590  gene are provided below the gene's name. Note, in the mapping from Ensembl (**Supplementary Data 3**) to
591  RefSeq annotations (required for use of ProteinPaint tool) a small number of variants annotations may have
592  changed or been lost, despite selecting the most similar representative transcript possible.
593

**References**

1.  Banerji S*, et al.* Sequence analysis of mutations and translocations across breast cancer subtypes. *Nature* **486**, 405-409 (2012).

2.  Stephens PJ*, et al.* The landscape of cancer genes and mutational processes in breast cancer. *Nature* **486**, 400-404 (2012).

3.  Cancer Genome Atlas N. Comprehensive molecular portraits of human breast tumours. *Nature* **490**, 61-70 (2012).

4.  Kennecke HF*, et al.* Late risk of relapse and mortality among postmenopausal women with estrogen responsive early breast cancer after 5 years of tamoxifen. *Annals of oncology : official journal of the European Society for Medical Oncology / ESMO* **18**, 45-51 (2007).

5.  Ellis MJ*, et al.* Whole-genome analysis informs breast cancer response to aromatase inhibition. *Nature* **486**, 353-360 (2012).

6.  Nielsen TO*, et al.* A comparison of PAM50 intrinsic subtyping with immunohistochemistry and clinical prognostic factors in tamoxifen-treated estrogen receptor-positive breast cancer. *Clinical cancer research : an official journal of the American Association for Cancer Research* **16**, 5222-5232 (2010).

7.  Chia SK*, et al.* A 50-gene intrinsic subtype classifier for prognosis and prediction of benefit from adjuvant tamoxifen. *Clinical cancer research : an official journal of the American Association for Cancer Research* **18**, 4465-4472 (2012).

8.  Samuels Y*, et al.* High frequency of mutations of the PIK3CA gene in human cancers. *Science* **304**, 554 (2004).

9.  Toy W*, et al.* ESR1 ligand-binding domain mutations in hormone-resistant breast cancer. *Nat Genet* **45**, 1439-1445 (2013).

10. Li S*, et al.* Endocrine-therapy-resistant ESR1 variants revealed by genomic characterization of breast-cancer-derived xenografts. *Cell Rep* **4**, 1116-1130 (2013).

11. Bose R*, et al.* Activating HER2 mutations in HER2 gene amplification negative breast cancer. *Cancer discovery* **3**, 224-237 (2013).

12. Amar D, Shamir R, Yekutieli D. Extracting replicable associations across multiple studies: Empirical Bayes algorithms for controlling the false discovery rate. *PLoS Computational Biology* **13**, e1005700 (2017).

13. Capanu M, Seshan VE. False discovery rates for rare variants from sequenced data. *Genetic epidemiology* **39**, 65-76 (2015).

14. Efron B. Size, Power and False Discovery Rates. *The Annals of Statistics* **35**, 1351-1377 (2007).

15. Gao J*, et al.* Integrative analysis of complex cancer genomics and clinical profiles using the cBioPortal. *Science signaling* **6**, pl1 (2013).

16. Curtis C*, et al.* The genomic and transcriptomic architecture of 2,000 breast tumours reveals novel subgroups. *Nature* **486**, 346-352 (2012).

647    17.    Pereira B*, et al.* The somatic mutation profiles of 2,433 breast cancers refines their genomic and
648            transcriptomic landscapes. *Nature communications* **7**, 11479 (2016).
649
650    18.    Endesfelder D*, et al.* Chromosomal instability selects gene copy-number variants encoding core
651            regulators of proliferation in ER+ breast cancer. *Cancer research* **74**, 4853-4863 (2014).
652
653    19.    McGranahan N, Burrell RA, Endesfelder D, Novelli MR, Swanton C. Cancer chromosomal instability:
654            therapeutic and diagnostic challenges. *EMBO reports* **13**, 528-538 (2012).
655
656    20.    Ma CX*, et al.* Neratinib Efficacy and Circulating Tumor DNA Detection of HER2 Mutations in HER2
657            Nonamplified Metastatic Breast Cancer. *Clinical cancer research : an official journal of the American
658            Association for Cancer Research* **23**, 5687-5695 (2017).
659
660    21.    Jaiswal BS*, et al.* Oncogenic ERBB3 mutations in human cancers. *Cancer cell* **23**, 603-617 (2013).
661
662    22.    Imielinski M*, et al.* Mapping the hallmarks of lung adenocarcinoma with massively parallel sequencing.
663            *Cell* **150**, 1107-1120 (2012).
664
665    23.    Shin DS*, et al.* Primary Resistance to PD-1 Blockade Mediated by JAK1/2 Mutations. *Cancer discovery*
666            **7**, 188-201 (2017).
667
668    24.    Zaretsky JM*, et al.* Mutations Associated with Acquired Resistance to PD-1 Blockade in Melanoma. *The
669            New England journal of medicine* **375**, 819-829 (2016).
670
671    25.    Fokkema IF, Taschner PE, Schaafsma GC, Celli J, Laros JF, den Dunnen JT. LOVD v.2.0: the next
672            generation in gene variant databases. *Human mutation* **32**, 557-563 (2011).
673
674    26.    Lukasik SM*, et al.* Altered affinity of CBF beta-SMMHC for Runx1 explains its role in leukemogenesis.
675            *Nature structural biology* **9**, 674-679 (2002).
676
677    27.    Mendes-Pereira AM*, et al.* Genome-wide functional screen identifies a compendium of genes affecting
678            sensitivity to tamoxifen. *Proceedings of the National Academy of Sciences of the United States of
679            America* **109**, 2730-2735 (2012).
680
681    28.    Yeh YC, Wu CC, Wang YK, Tang MJ. DDR1 triggers epithelial cell differentiation by promoting cell
682            adhesion through stabilization of E-cadherin. *Molecular biology of the cell* **22**, 940-953 (2011).
683
684    29.    Rudd ML*, et al.* Mutational analysis of the tyrosine kinome in serous and clear cell endometrial cancer
685            uncovers rare somatic mutations in TNK2 and DDR1. *BMC cancer* **14**, 884 (2014).
686
687    30.    Loriaux MM*, et al.* High-throughput sequence analysis of the tyrosine kinome in acute myeloid
688            leukemia. *Blood* **111**, 4788-4796 (2008).
689
690    31.    Ding L*, et al.* Somatic mutations affect key pathways in lung adenocarcinoma. *Nature* **455**, 1069-1075
691            (2008).
692
693    32.    Vogel WF, Aszodi A, Alves F, Pawson T. Discoidin domain receptor 1 tyrosine kinase has an essential
694            role in mammary gland development. *Molecular and cellular biology* **21**, 2906-2917 (2001).
695
696    33.    Peng Y*, et al.* Deficiency in the catalytic subunit of DNA-dependent protein kinase causes down-
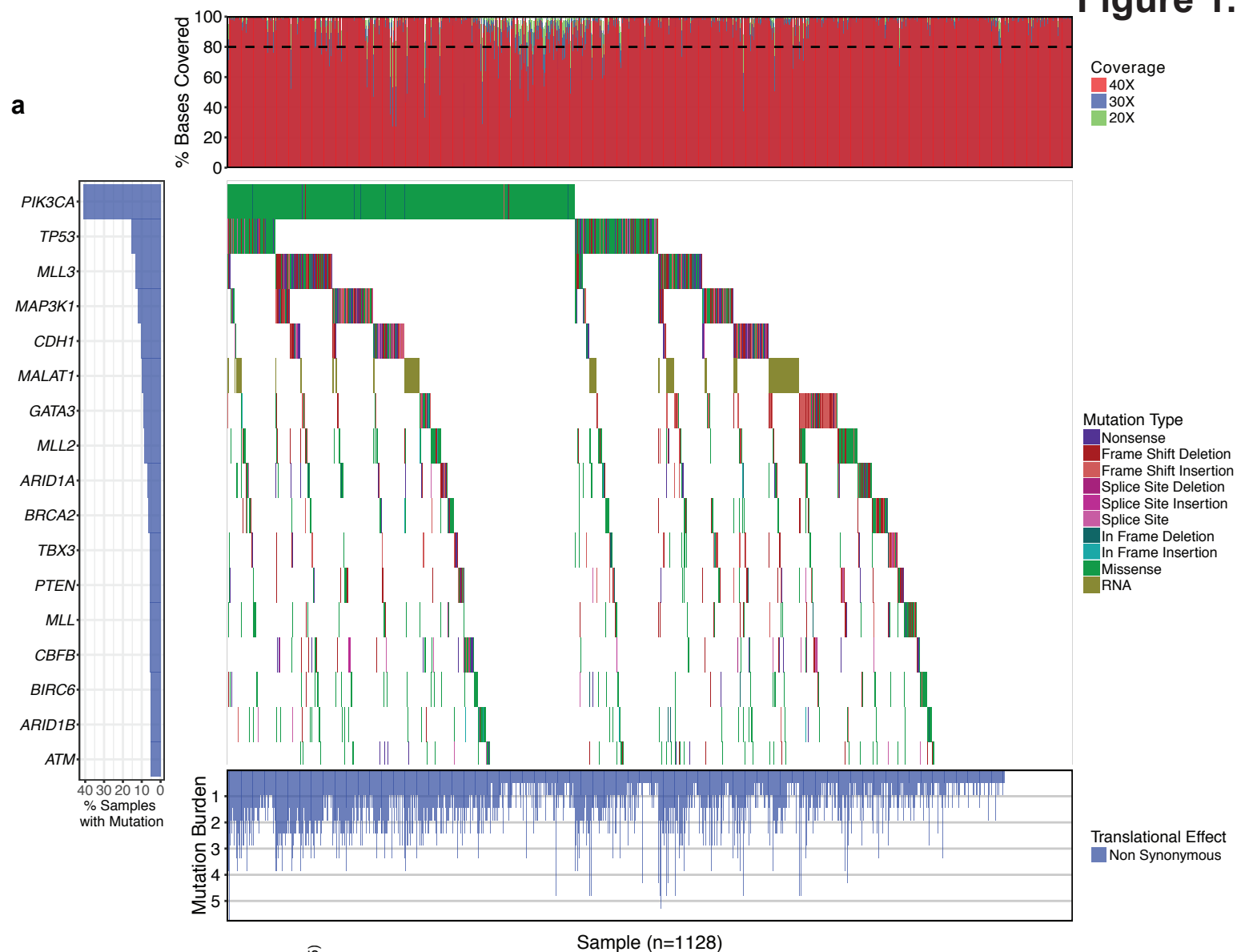697            regulation of ATM. *Cancer research* **65**, 1670-1677 (2005).
698

34. Haricharan S*, et al.* Loss of MutL Disrupts CHK2-Dependent Cell-Cycle Control through CDK4/6 to Promote Intrinsic Endocrine Therapy Resistance in Primary Breast Cancer. *Cancer discovery* **7**, 1168-1183 (2017).

35. Cheang MC*, et al.* Ki67 index, HER2 status, and prognosis of patients with luminal B breast cancer. *Journal of the National Cancer Institute* **101**, 736-750 (2009).

36. Liu S*, et al.* Prognostic significance of FOXP3+ tumor-infiltrating lymphocytes in breast cancer depends on estrogen receptor and human epidermal growth factor receptor-2 expression status and concurrent cytotoxic T-cell infiltration. *Breast cancer research : BCR* **16**, 432 (2014).

37. Parker JS*, et al.* Supervised risk predictor of breast cancer based on intrinsic subtypes. *Journal of clinical oncology : official journal of the American Society of Clinical Oncology* **27**, 1160-1167 (2009).

38. Kan Z*, et al.* Diverse somatic mutation patterns and pathway alterations in human cancers. *Nature* **466**, 869-873 (2010).

39. Shah SP*, et al.* The clonal and mutational evolution spectrum of primary triple-negative breast cancers. *Nature* **486**, 395-399 (2012).

40. Griffith M*, et al.* DGIdb: mining the druggable genome. *Nature methods* **10**, 1209-1210 (2013).

41. Chanock SJ*, et al.* Somatic sequence alterations in twenty-one genes selected by expression profile analysis of breast carcinomas. *Breast cancer research : BCR* **9**, R5 (2007).

42. Griffith M*, et al.* Genome Modeling System: A Knowledge Management Platform for Genomics. *PLoS Comput Biol* **11**, e1004274 (2015).

43. Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* **25**, 1754-1760 (2009).

44. Li H*, et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078-2079 (2009).

45. Lek M*, et al.* Analysis of protein-coding genetic variation in 60,706 humans. *Nature* **536**, 285-291 (2016).

46. Genomes Project C*, et al.* A map of human genome variation from population-scale sequencing. *Nature* **467**, 1061-1073 (2010).

47. Fu W*, et al.* Analysis of 6,515 exomes reveals the recent origin of most human protein-coding variants. *Nature* **493**, 216-220 (2013).

48. Cancer Genome Atlas Research N. Genomic and epigenomic landscapes of adult de novo acute myeloid leukemia. *The New England journal of medicine* **368**, 2059-2074 (2013).

49. Krysiak K*, et al.* A genomic analysis of Philadelphia chromosome-negative AML arising in patients with CML. *Blood cancer journal* **6**, e413 (2016).

50. Ainscough BJ*, et al.* DoCM: a database of curated mutations in cancer. *Nature methods* **13**, 806-807 (2016).

51. Zhou X*, et al.* Exploring genomic alteration in pediatric cancer using ProteinPaint. *Nat Genet* **48**, 4-6 (2016).
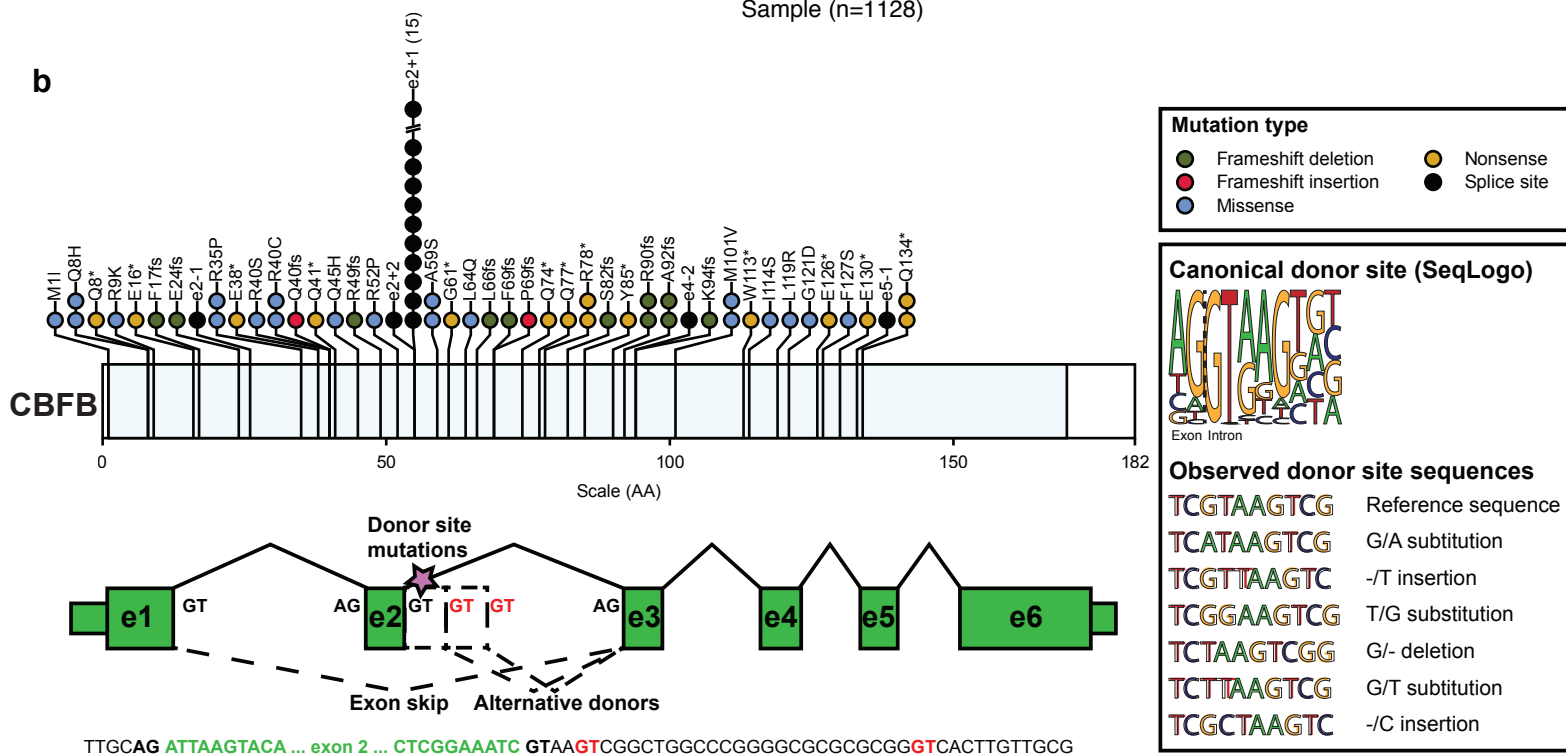
753    52.    Skidmore ZL*, et al.* GenVisR: Genomic Visualizations in R. *Bioinformatics* **32**, 3012-3014 (2016).
754
755

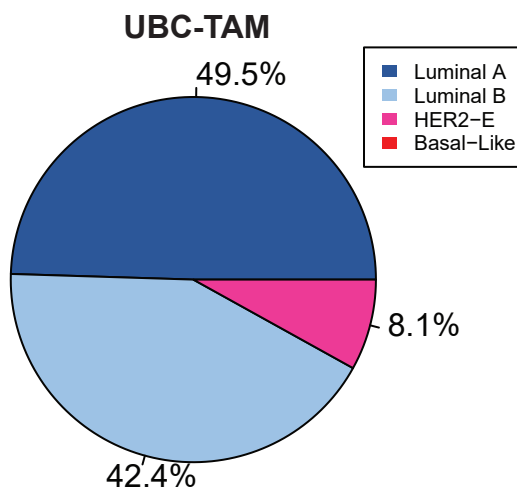**Figure 1.**
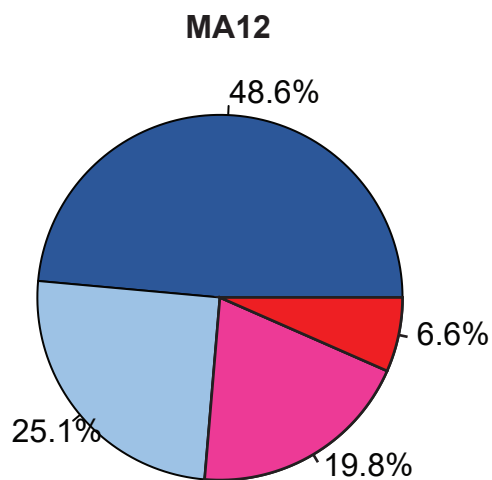
Figure 2.

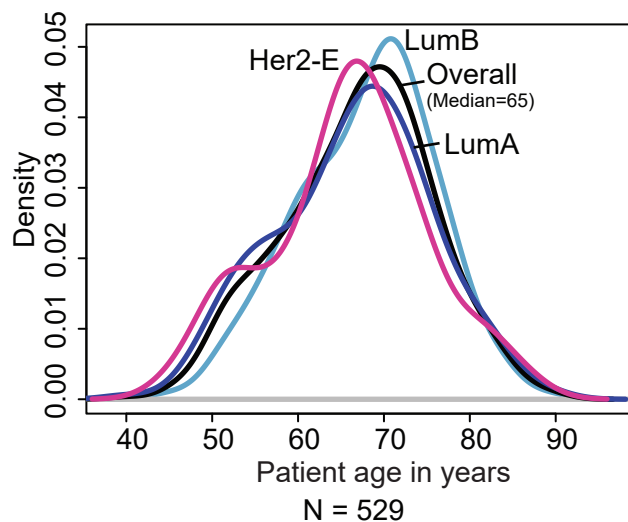**a**

**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    Overall (Median=43)    LumB    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 Cohort**
- 947 Tam-treated patients.
- 645 with >50ng DNA.
- 632 ER+, tumor-only.

Gene Selection:
- 83 genes.
- 8 studies' recurrently mutated genes

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

**Data Analysis**
- Mutation Landscape.
- Survival Analysis.
- Mutational 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 | #Cases | p-value |
|------|---------|--------|---------|
| ARID1B | nonsilent | 30 | 0.1410 / 0.0310 |
| DDR1 | nonsilent | 8 | 2.42E-05 / 0.0047 |
| ERBB3 | nonsilent | 17 | 0.0606 / 0.0841 |
| MAP3K1 | nonsilent | 85 | 0.0033 / 0.0245 |
| NF1 | FS/NS | 7 | 0.1256 / 0.0167 |
| PIK3CA | nonsilent | 285 | 0.0260 / 0.0349 |
| PRKDC | nonsilent | 18 | 0.0380 / 0.0179 |
| TP53 | FS/NS | 19 | 0.0288 / 0.0127 |
| TP53 | nonsilent | 80 | 0.0216 / 0.0899 |
| XBP1 | nonsilent | 14 | 0.0683 / 0.0593 |

BCSS — RFS

**Univariate Hazard Ratio**

## c



| Gene | Variant | #Cases | p-value |
|------|---------|--------|---------|
| ARID1B | nonsilent | 30 | 0.3145 |
| DDR1 | nonsilent | 8 | 0.0000 |
| ERBB3 | nonsilent | 17 | 0.0687 |
| GATA3 | FS/NS | 34 | 0.0364 |
| MAP3K1 | nonsilent | 85 | 0.0014 |
| MET | nonsilent | 8 | 0.0859 |
| NF1 | FS/NS | 7 | 0.2322 |
| PIK3CA | nonsilent | 285 | 0.0348 |
| PRKDC | nonsilent | 18 | 0.4917 |
| TP53 | FS/NS | 19 | 0.0333 |
| TP53 | nonsilent | 80 | 0.0389 |
| Tumor Grade | clinical | | 0.0026 |
| Node positivity | clinical | | 0.0000 |
| Tumor Size >5cm | clinical | | 0.0000 |

**Multivariate Hazard Ratio**

**Figure 4.**

**a**



**b**

**Figure 5.**

**a**

UBC-TAM Candidates
*(UVA p<0.05)*

ARID1B    PIK3CA
DDR1      PRKDC
ERBB3     TP53
MAP3K1    XBP1
NF1

METABRIC dataset

- Mutations extracted from cbioportal (Pereira, Nat Comm 2016)
- Clinical annotations linked using Oncomine Curtis Breast 2 dataset (Curtis, Nature 2012)
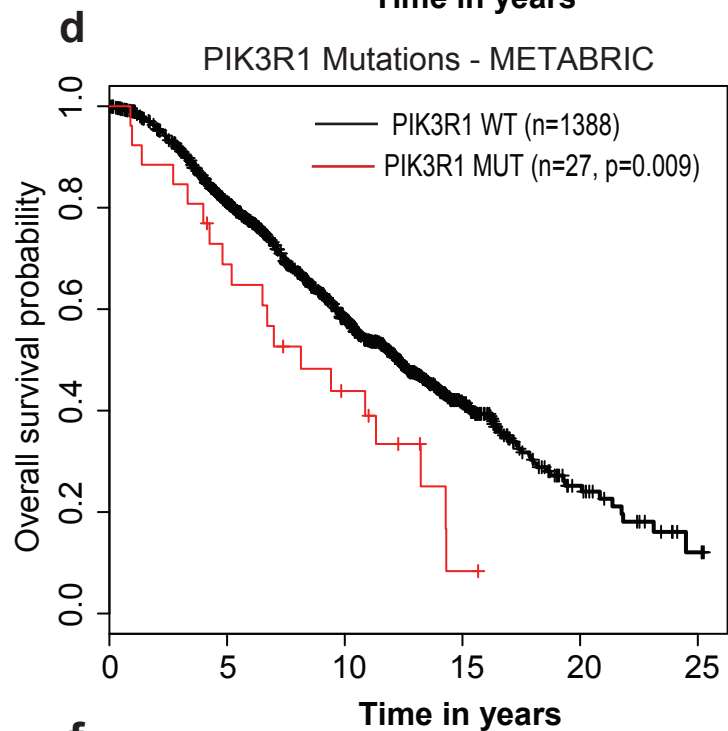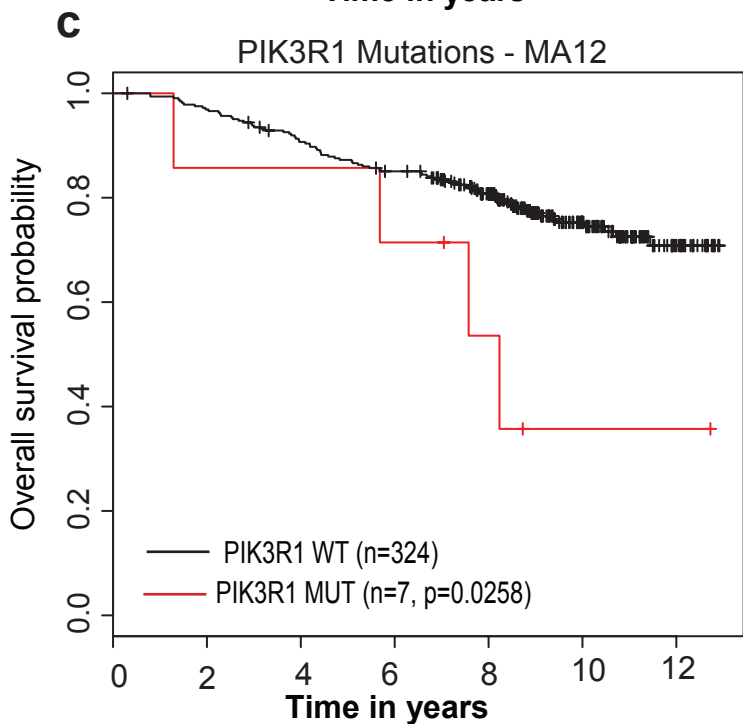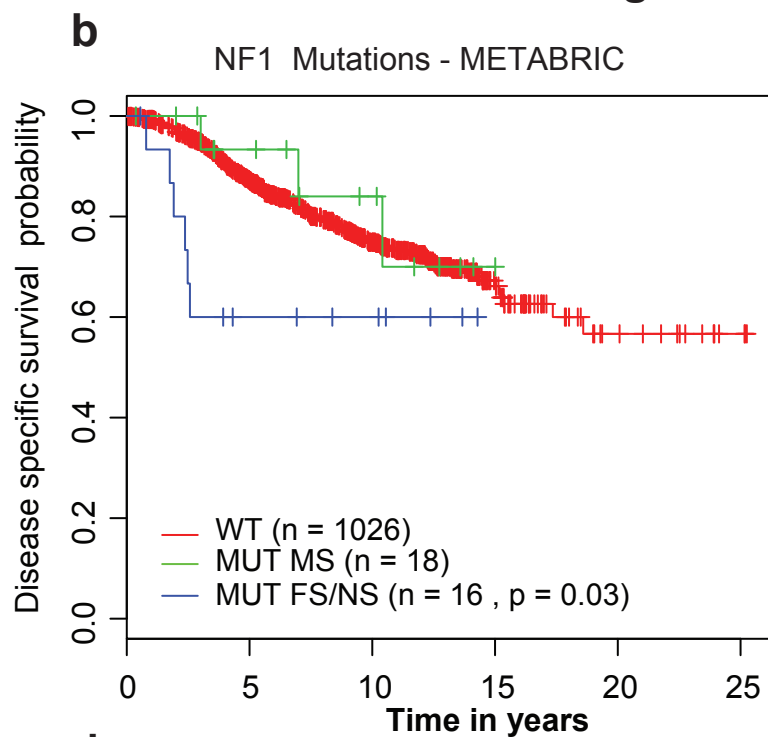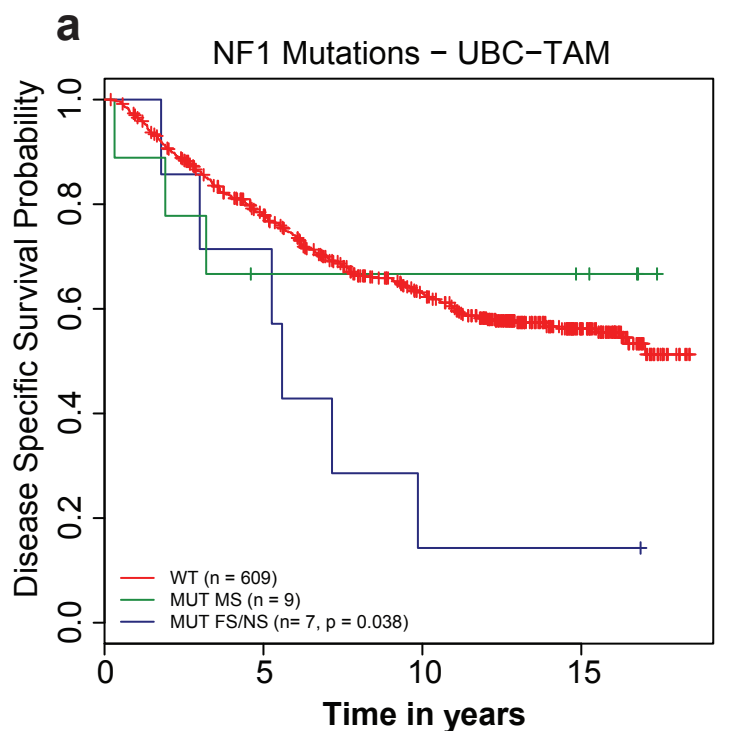
Univariate analysis MVA with Grade, Tumor Size, Node Status

Validation

Metabric Cohort (2369) → ER+ (1415) → Removed samples present in UBC-TAM (1148) → Disease-specific survival information available (1060)

**b**

| Gene | Variant | #cases/CNV | P-value |
|------|---------|-----------|---------|
| ARID1B | nonsilent | 25 | 0.3102 / 0.269 |
| ERBB3 | nonsilent | 30 | 0.8288 / 0.852 |
| MAP3K1 | nonsilent | 130 | 0.009239 / 0.0174 |
| NF1 | FS/NS | 16 | 0.03 / 0.011 |
| PIK3CA | nonsilent | 500 | 0.6438 / 0.903 |
| TP53 | nonsilent | 248 | 2.17e−12 / 5.14e−09 |
| Tumor Grade | Clinical | | 2.21e−05 |
| Node Positivity | Clinical | | 4.71e−09 |
| Tumor Size >5cm | Clinical | | 7.2e−07 |

Univariate BCSS    Multivariate BCSS

No_CNV
Deep_del    Shallow_del
Amp    Gain

Hazard Ratio

**Figure 6.**

**Figure 7.**

Figure 8.