# Supplemental Figures for Shah et al: The clonal and mutational evolution spectrum of primary triple negative breast cancers

Figure 1: Distribution of sequence coverage of targeted exons for 54 tumour and normal exome capture libraries. Median of median coverage over targeted exons over all libraries was 60.5x.

Figure 2: Distribution of the number of mutations per sample as it relates to cellularity according to pathologist review

Figure 3: Somatic copy number aberration (CNA) landscape of triple negative breast cancers. **(a)** The somatic CNA landscape of Basal (top) and Other (bottom) triple negative cancers showing significant segmental CNA instability across the genome. The frequency alteration of each gene by Ensembl54 annotation is shown where deletions (blue) are reflected as negative frequencies and gains (red) are shown as positive frequencies. **(b)** Comparison of genomic instability between Basal and non-Basal cancers for each of the alteration states. Comparisons are shown as boxplots of the distribution of proportion of the genome altered for each CNA state over the population of 80 tumours that were classifiable by PAM50. All states except HLAMPs showed a statistically higher proportion of alteration in Basals than Others. **(c)**Proportion of homozygous somatic SNV mutations overlapping segmental CNAs, segregated by state. **(d)** RNASeq derived gene fusions supported by segmental CNA breakpoints in Basal (left) and Other (right) cancers depicted as circos diagrams. Chromosomes are arrayed around the circle and arcs between points on the circle represent gene fusions predicted by defuse [1].

Figure 4: The somatic CNA landscape of Basal (top) and Other (bottom) triple negative cancers from an external cohort of 122 TNBC tumours. The frequency alteration of each gene by Ensembl54 annotation is shown where deletions (blue) are reflected as negative frequencies and gains (red) are shown as positive frequencies. Patterns of alteration were consistent with the cohort of 104 TNBCs featured in this study.

Figure 5: Predicted intragenic homozygous deletions in PARK2 on chr6q26 inferred from Affymetrix SNP 6.0 arrays. Predictions for 6 cases are shown where vertical grey bars denote the boundaries of PARK2. Orange box on the ideogram (bottom right) indicates the genomic coordinates of the view. Each datapoint corresponds to the log ratio of representing total copy number. Bright green points indicate homozygous deletions inferred by segmentation by HMM-Dosage (Supplemental methods). Dark green points are predicted hemizygous regions and blue points are predicted neutral regions.

{ sfig:park2˙homd }

Figure 6: Examples of four (multi-page PDF) somatic mutation-associated alternative splicing events indicated on the distribution of canonical and alternative exon junction usage. Specific cases that have a donor/acceptor dinucleotide splice site mutation are indicated as red dots. Distributions are shown for splice site mutation events in TP53, PIK3R1, LAMA2 and LAMB1. In all examples, cases harbouring a mutation in the splice sites show outlying alternative junction usage.

{ sfig:splice˙site˙examples }

Figure 7: List of deleterious variations occurring within RB predicted for the triple negative breast cancer set. The coordinates of the variations and of functional TFBS predicted with the reference sequence (without variation) are provided in the hg18 assembly. The most proximal gene is based on Ensembl annotation (release 63). In the two columns with sequences, the red nucleotides indicate a variation position. The last column highlights the variant position within the TFBS logo (red rectangle). The data are ranked by the ratio of the relative scores. The lower this value, the greater the predicted alteration to the recognition of the sequence by the TF. The rows with bold text reflect the most damaging impact predicted (ratio of relative score lower than 0.9 - see Supplemental methods).

{ sfig:tfbs˙mutation }

Figure 8: Clustering of extreme CNA event (HOMD, HLAMP) into gene modules and functional pathways by Reactome analysis. Nodes in the network are genes affected by a high level amplification or homozygous deletion. Size of the node is proportional to the frequency of abberation in the population of 104 tumours. Edges between nodes represent protein interactions as described in the Reactome FI database [2]

{ sfig:CNA˙modules˙pathway

Figure 9: Ab-initio clustering of the mutational profiles of cases and pathways. Pathway-based features were calculated based on the frequency of mutation in each pathway for a given case. Conservatively, we included 48 pathways showing enrichment from the entire network for further analysis (Table S14) containing 40 somatically mutated genes. For each tumour, we counted the number of times a gene in a given pathway was found to be mutated, producing a tumour-pathway matrix. Hierarchical clustering reveals a population structure where distinct groupings of cases by their pathway-based mutational profiles are evident. Annotations of the number of mutations in each case (scaled in blue), the PAM50 subtype and TP53, Rb1 and PTEN mutation status are shown across the top of the heatmap. Rows in the matrix correspond to pathways shown in Figure 3b. The relationship between genes and pathways in this analysis is such that a small number of highly connected gene nodes (drivers) will implicate many pathways, however the strength of pathway association is also modified by the presence of additional mutated genes, shown as a grayscale for each tumour-pathway data point.

{ sfig:networkclustering }

Figure 10: Probabilistic graphical model of a Dirichlet process for clonal frequency estimation. Left the random variables in the model are shown as nodes in a graph where shade nodes are observed quantities and unshaded nodes are unobserved quantities we wish to estimate. The model inputs the total number of reference and variant reads at the mutation of interest, the observed copy number at the genomic position of each mutation $i$ and then infers the genotype of the mutation ($G_i$), the clonal frequency ($\phi_i$) and the number of variant reads belonging to each clone ($A_{var}^i$, $D_{var}^i$). The conditional probability distributions for the model are shown (right).

{ sfig:pyclone }

Figure 11: Clonal frequency analysis for all tumours with 10 more mutations in the TNBC cohort (multipage PDF). This figure shows on alternating pages, the clustering of the co-occurrence matrix for mutations to infer the assignment to clonal clusters and on the subsequent page the full inferred posterior distributions of the clonal frequency estimates as output from 2500 runs of the MCMC sampling of the Dirichlet process model. The mutations in the distribution plots are first ordered by the average clonal frequencies over each group of mutations. The mode of the clonal frequency estimate for each mutation is shown with a grey dot. { sfig:omnibus`evolution }

Figure 12: Distribution of clonal frequencies across all mutations (adjusted for predicted normal contamination). Specific clonal frequencies are shown for TP53 mutations and PIK3CA mutations and reveals that the distribution of TP53 mutation is incompatible with the notion that it is always the founder mutation (and therefore present in all cells). Instead the distribution is skewed in favour of TP53 residing in a dominant proportion of cells, but exhibits a tail suggesting that in some tumours it occurs later in the evolution of the cancer. { sfig:tp53clonality }

# References

[1] McPherson, A. *et al.* deFuse: An Algorithm for Gene Fusion Discovery in Tumor RNA-Seq Data. *PLoS Comput Biol* **7** (2011).

[2] Wu, G., Feng, X. & Stein, L. A human functional protein interaction network and its application to cancer data analysis. *Genome Biology* **11**, R53 (2010).

Figure S1

Figure S2

(a) Basal / Other — genome-wide copy number frequency plots

Figure S3

(b) Total, HOMD, HETD, GAIN, AMP, HLAMP boxplots (Basal vs Other)
- Total: p-value=1e−04
- HOMD: p-value=0.0104
- HETD: p-value=0.0107
- GAIN: p-value=8e−04
- AMP: p-value=6e−04
- HLAMP: p-value=0.3244

(c) SNP6 derived copy number states, pval=3.58e−11
HOMD, HETD, NEUT, GAIN, AMP, HLAMP
ALL, HOM, HET
Proportion of mutations

(d) Basal

(e) Other

Figure S4

Figure S5

Figure 6
4 pages

**SA089, mutation chr17:7517653 gene TP53**



Canonical Junction

Alternative Junction  1

Alternative Junction  2

**SA208, mutation chr5:67627003 gene PIK3R1**

**SA097, mutation chr6:129517341 gene LAMA2**

Canonical Junction

Alternative Junction

Normalized junction read count

chr6:129511934−129517343

chr6:129511934−129517364

SA097

| ID | Hg18 chromosome | Hg18 coordinate | TFBS coordinates | Closest gene Ensembl ID and gene name | Reference sequence | Reference relative score | Alternative sequence | Alternative relative score matrix | Ratio of relative score | | Location of the variation in the TFBS |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | | Figure S7 |
| SA073 | 1 | 215871678 | 215871667 215871681 | ENSG00000162814 SPATA17 | CTTTTTTTCTT**C**GAC | 80.7 | CTTTTTTTCTT**T**GAC | 70.9 | 0.88 | C>T |  |
| SA236 | 1 | 24401404 | 24401397 24401411 | ENSG00000230023 RP11-10N16.2 | TTCTTGT**C**CAGGCCC | 83.0 | TTCTTGT**G**CAGGCCC | 73.2 | 0.88 | C>G |  |
| SA065 | 21 | 34656187 | 34656178 34656192 | ENSG00000159197 KCNE2 | CCCTGTTCC**C**ACCAG | 88.2 | CCCTGTTCC**G**ACCAG | 78.6 | 0.89 | C>G |  |
| SA065 | 19 | 48793075 | 48793064 48793078 | ENSG00000226763 SRRM5 | TGCACGCTCCT**C**TCC | 82.1 | TGCACGCTCCT**G**TCC | 73.9 | 0.90 | C>G |  |
| SA065 | 19 | 15390963 | 15390957 15390971 | ENSG00000246896 AC020911.1 | TGAGAG**C**CCCGCCTT | 82.9 | TGAGAG**G**CCCGCCTT | 75.2 | 0.91 | C>G |  |
| SA237 | 5 | 97673262 | 97673262 97673276 | ENSG00000176183 CH17-12M21.1 | **G**TCTTTTCCTTTCTC | 81.1 | **A**TCTTTTCCTTTCTC | 79.1 | 0.98 | G>A |  |

Figure S8

Figure S9

$$\phi^i \quad \sim \quad DP(\alpha, G_0)$$
$$\quad = \quad \text{Clonal frequency}$$
$$D_{var}^i | \phi^i, d^i \quad \sim \quad Binomial(d_{var}^i | \phi^i, d^i)$$
$$\quad = \quad \text{Number of reads sampled from clones}$$
$$A_{var}^i | G^i, \mu_{var}, D_{var}^i \quad \sim \quad Binomial(a_{var}^i | \mu_{var:cg}, d_{var}^i)$$
$$\quad = \quad \text{Number of matches to ref from clones}$$
$$G^i | C^i \quad \sim \quad Discrete\,Uniform(g | 1, c)$$
$$\quad = \quad \text{Genotype of clone}$$
$$a^i | A_{var}^i, d^i, D_{var}^i, \mu_{ref} \quad \sim \quad Binomial(a^i - a_{var}^i | \mu_{ref}, d^i - d_{var}^i)$$
$$\quad = \quad \text{Number of matches to ref}$$

SA018 n= 33

Figure S11
108 pages

**SA018 n= 33**

Clonal frequency

SA029 n= 11

**SA029 n= 11**

Clonal frequency

SA030 n= 22

SA030 n= 22

Clonal frequency

SA031 n= 60

SA031 n= 60

Clonal frequency

SA051 n= 34

**SA051 n= 34**

Clonal frequency

SA052 n= 25

SA052 n= 25

SA054 n= 70

SA054 n= 70

Clonal frequency

SA056 n= 12

Clonal frequency

SA067 n= 21

**SA067 n= 21**

Clonal frequency

SA069 n= 54

SA069 n= 54

Clonal frequency

SA071 n= 96

**SA071 n= 96**

Clonal frequency

**SA072 n= 22**
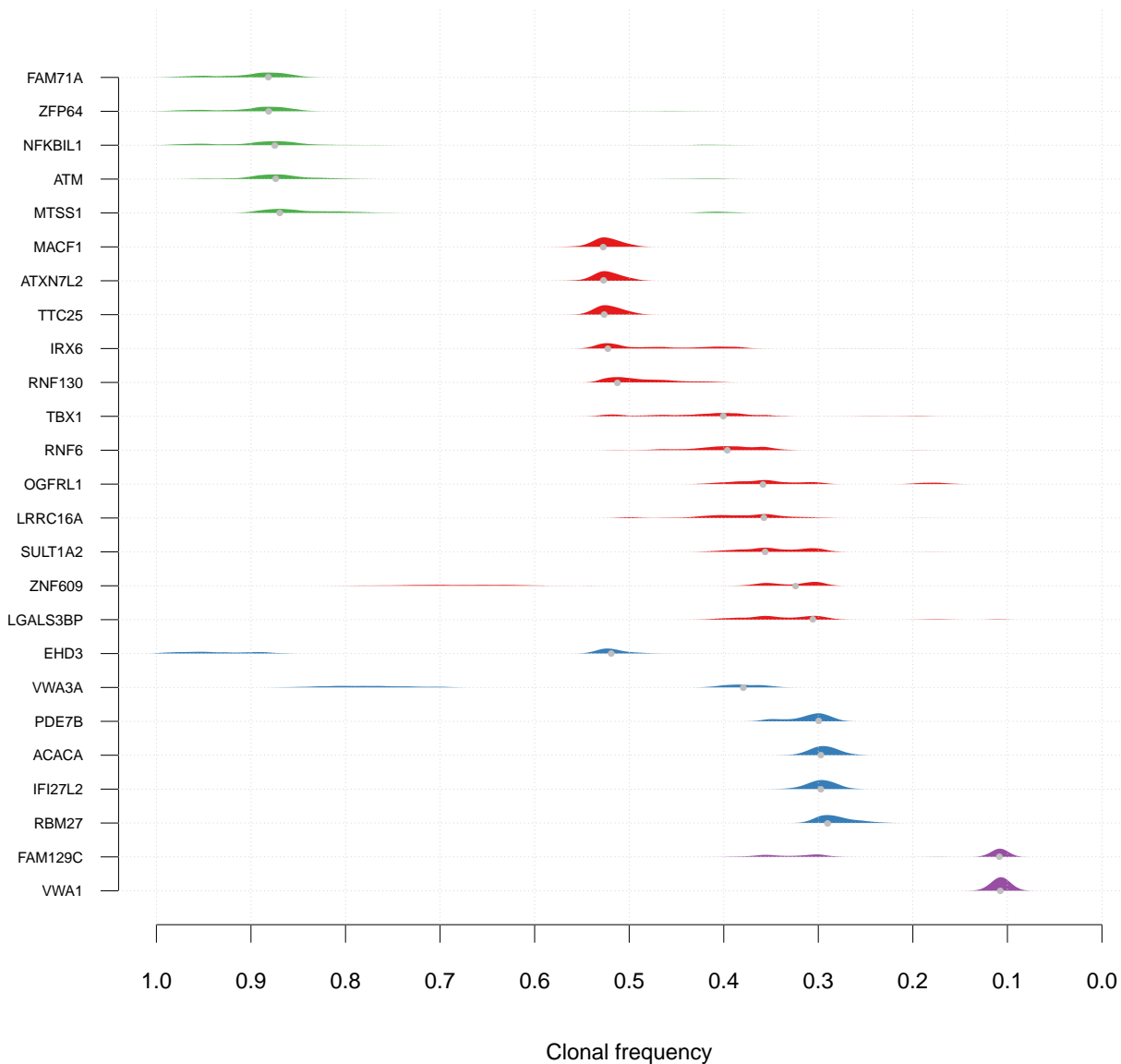
Clonal frequency

SA073 n= 14

SA073 n= 14

**SA074 n= 24**

Clonal frequency

SA075 n= 34

SA075 n= 34

SA076 n= 14

SA076 n= 14

Clonal frequency

SA077 n= 77

**SA077 n= 77**

Clonal frequency

SA083 n= 50

SA083 n= 50

SA084 n= 65

**SA084 n= 65**

Clonal frequency

SA089 n= 49

SA089 n= 49

Clonal frequency

SA090 n= 36

**SA090 n= 36**

Clonal frequency

SA092 n= 66

SA092 n= 66

SA093 n= 17

SA093 n= 17

Clonal frequency

SA094 n= 26

**SA094 n= 26**

Clonal frequency

**SA096 n= 10**

Clonal frequency

SA097 n= 46

**SA097 n= 46**

Clonal frequency

**SA098 n= 22**

Clonal frequency

SA101 n= 18

Clonal frequency

SA102 n= 24

Clonal frequency

SA103 n= 15

**SA103 n= 15**

Clonal frequency

SA106 n= 132

**SA106 n= 132**

Clonal frequency

SA208 n= 27

**SA208 n= 27**

Clonal frequency

SA210 n= 56

SA210 n= 56

Clonal frequency

SA212 n= 48

**SA212 n= 48**

Clonal frequency

SA213 n= 44

**SA213 n= 44**

Clonal frequency

**SA214 n= 195**

SA214 n= 195

Clonal frequency

SA215 n= 12

**SA215 n= 12**

Clonal frequency

SA216 n= 16

**SA216 n= 16**

Clonal frequency

SA217 n= 14

SA217 n= 14

Clonal frequency

SA218 n= 112

SA218 n= 112

Clonal frequency

SA219 n= 36

**SA219 n= 36**

Clonal frequency

SA222 n= 15

**SA222 n= 15**

Clonal frequency

**SA223 n= 16**

SA223 n= 16

Clonal frequency

SA225 n= 60

SA225 n= 60

Clonal frequency

SA226 n= 20

SA226 n= 20

Clonal frequency

SA227 n= 16

**SA227 n= 16**

SA228 n= 16

**SA228 n= 16**

Clonal frequency

SA229 n= 33

SA229 n= 33

Clonal frequency

SA230 n= 10

SA231 n= 26

Clonal frequency

SA234 n= 21

**SA234 n= 21**

Clonal frequency

**SA235 n= 17**

Clonal frequency

SA236_n= 45

SA236 n= 45

Clonal frequency

SA237 n= 55

Clonal frequency

Figure S12

# Supplemental Methods for Shah et al: The clonal and mutational evolution spectrum of primary triple negative breast cancers

## 1  Biospecimen collection and ethical consent

We assembled a collection of 149 clinically annotated primary fresh frozen triple negative breast cancer specimens ($n$ = 149) each with a source of matched non malignant tissue for extraction of germline representative DNA. Tumour specimens were obtained from three tumour banks (BCCA Vancouver, breast tumour tissue repository; Alberta CBCF Breast Tumour Bank Edmonton, Cambridge UK, Addenbrooke's Hospital breast tumour bank) each with local REB/IRB approval for genomic studies of nucleic acids from breast cancer patients. This project was conducted under local BCCA REB/IRB projects H06-00289, H08-1230, H06-3199. The source of germline DNA was from peripheral blood lymphocytes in all but 4 cases. In these 4 cases histologically normal adjacent breast tissue was used. Initial case selection was based on clinical immunohistochemistry to define primary triple negative breast cancers obtained from surgical specimens, prior to the initiation of any chemotherapy or radiotherapy. Tumours typed as ER-, HER2- and PR- were initially selected for further review and re-validation of the IHC.

## 2  Histopathological review

Tissue sections were subject to expert histopathological review (GT) to assess the presence of invasive tumour, pre-malignant or benign changes, lymphocytic infiltration, necrosis and tumour cellularity. Tumour cellularity was scored visually in a semiquantitative fashion on sections taken from the cryosectioning runs used to isolate nucleic acids from each tumour. Cellularity values were binned such that 'low cellularity' corresponds to samples with <40% malignant cells, 'moderate cellularity' corresponds to 40% - 70% malignant cells, and samples with >70% malignant cells were considered to have 'high cellularity'. All but one sample classified as low cellularity were excluded from further analysis. The ER-, PR- HER2- immunophenotype was reconfirmed on sections or TMA cores from the cases included for analysis and additionally CK5/6 and EGFR were assessed by IHC. Subsequently SNP6.0 copy number analysis was also used to confirm the absence of HER2 amplification in each case.

## 3  Sequence data generation

Paired end RNASeq libraries were generated as described in Wiegand et al [1]. Exome capture was achieved through solution hybrid selection with the Human All Exon kit SureSelect Target Enrichment System (Agilent) version 1 for Illumina Genome Analyzer paired-end sequencing [2] and libraries were prepared as described in [3]. SOLiD whole genome shotgun libraries for 15 tumour/normal pairs were generated as previously described [4]. Details of per case coverage for each sequencing method are presented in Table S2. Targeted sequencing interrogating the exons and splice sites of 29 genes in 159 additional breast cancers (82 ER+ and 77 ER-ve, tumour and matched normal) was undertaken using exon-specific PCR and cDNA probe-based capture. The amplified/captured products were sequenced on the Illumina HiSeq 2000 platform.

# 4 Exome sequence data analysis and mutation validation

**Sequence alignment and mutation calling**   Sequence reads were aligned to the human reference sequence (NCBI build 36, hg18) and a compacted reference composed of the targeted exons. Alignments were performed using Maq v0.7.1 [5] using default parameters. Uniquely aligned reads were kept for downstream analysis. Somatic SNVs were called using JointSNVMix (`http://compbio.bccrc.ca/software/jointsnvmix/`) [6] combined with mutationSeq [7], from the resulting tumour and normal pair of BAM files for each case. We used a marginal posterior probability threshold of 0.2 to nominate somatic SNVs, aiming to err on the side of sensitivity and minimize false negatives. All non-synonymous mutations as profiled by MutationAssessor [8] (`http://mutationassessor.org`) were carried forward for validation by deep amplicon sequencing. Indels were predicted using samtools v0.1.7 pileup command from BAM files. Indels present in at least 10% of reads covering the position of interest, predicted to affect a coding sequence, and for which there were 0 reads containing the same indel in the corresponding normal DNA library were carried forward for validation as candidate somatic indels.

## 4.1   Validation of mutations

**Revalidation sequencing**   Somatic mutation predictions were subjected to rigorous validation by targeted deep amplicon sequencing in the tumour and normal DNA templates. Genomic DNA was prepared as previously described [9] . Automated primer design was performed using Primer3 (Rozen and Skaletsky, 1998) and custom scripting. Primer pairs were designed to place the variant position within 75bps of either end of the amplicon and to be between 50-300bp in length. Primer pairs were independently validated by in silco PCR followed by BLAT against the human genome to ensure that the correct target was generated and that the resulting amplicon was unique within the genome. DNA primers were synthesized in 96-well plates at a 25nmol scale with standard desalting (IDT Coralville, IA USA). Polymerase cycling reactions were set up in 96-well plates and comprised of 0.5 $\mu$M forward primer, 0.5 $\mu$M reverse primer, 1-2 ng of gDNA template, 5X Phusion HF Buffer, 0.2 $\mu$M dNTPs, 3% DMSO, and 0.4 units of Phusion DNA polymerase (NEB, Ipswich, MA, USA). Reaction plates were cycled on a MJR Peltier Thermocycler (model PTC-225) with cycling conditions of a denaturation step at $98\,°C$ for 30 sec, followed by 35 cycles of [$98\,°C$ for 10 sec, $69\,°C$ for 15 sec, $72\,°C$ for 15 sec] and a final extension step at $72\,°C$ for 10 min. PCR reactions were visualized on 3% agarose (NuSieve) gels for 2hrs at 170V to assess PCR success. Successful reactions were manually pooled (4ul per well) by template and subjected to Illumina library construction using a modified paired-end protocol (Illumina, Hayward, USA). This involved A-tailing of the amplicons and ligation to Illumina PE adapters. Adapter-ligated products were purified on Qiaquick spin columns (Qiagen, Valencia, CA, USA) and PCR-amplified using Phusion DNA polymerase (NEB, Ipswich, MA, USA) in 10 cycles using PE primer 1.0 (Illumina) and a custom multiplexing PCR Primer [5-CAAGCAGAAGACGGCATACGAGAT NNNNNNCGGTCTCGGCATTCCTGC TGAACCGCTCTTCCGATCT-3] where *NNNNNN* was replaced with unique fault tolerant hexamer barcodes for each template. PCR products of the desired size range were purified away from adapter ligation artifacts using 8% PAGE gels, pooled and DNA quality was assessed and quantified using an Agilent DNA 1000 series II assay (Agilent, Santa Clara CA, USA) and Nanodrop 7500 spectrophotometer (Nanodrop, Wilmington, DE, USA) and subsequently diluted to 10nM. The final concentration was confirmed using a Quant-iT dsDNA HS assay kit and Qubit fluorometer (Invitrogen, Carlsbad, CA, USA). For sequencing, clusters were generated on the Illumina cluster station using v4 cluster reagents and paired-end 75bp reads generated using v4 sequencing reagents on the Illumina GA$_{ii}$ platform following the manufacturer's instructions. Between the paired 75bp reads a third 7 base pair read was performed using the following custom sequencing primer [5- GATCGGAAGAGCGGTTCAGCAGGAATGCCGAGACCG] to sequence the hexamer barcode. Image analysis, base-calling and error calibration was performed using v1.60 of Illumina's Genome analysis pipeline.

**Revalidation data analysis**   PCR primers were designed to revalidate candidate somatic mutations (SNVs and indels) on whole genome amplified material for both tumour and normal DNA. Resulting amplicons were indexed, pooled and sequenced deeply using 75bp paired end libraries as described above. Reads were separated out by 6bp index using a Hamming distance function with a maximum mismatch tolerance of 1,

and aligned (Maq v0.7.1) to a reference database containing only the targeted loci for each case. (Note, we compared results to BWA [10] and there were no observable differences in the validated list of mutations). For each targeted sequence, we inferred the presence/absence of the targeted variants using a Binomial exact test as previously described [9]. In addition, we imposed the following criteria in order to classify a 'validated' somatic mutation:

- both tumour and normal data had a minimum of 50 reads covering the targeted position

- the Binomial exact test result (Benjamini Hochberg adjusted p-value) for the tumour $< 0.01$

- the Binomial exact test result (Benjamini Hochberg adjusted p-value) for the normal $>= 0.01$

- the proportion of reads indicating the variant in the tumour $\geq 5\%$

For indels, we required:

- a minimum of 10% of all reads aligned to the position containing the indel

- at least 3 reads containing the indel

- a minimum local realignment score of 300

One validation library was single end. For this library, we aligned the reads using GSNAP [11] - a gapped aligner capable of detecting indels in single end reads.

## 5   RNASeq data analysis and validation

**Sequence alignment and SNV/indel calling for RNASeq**   SNVs for RNASeq were analysed as previously described [9] using SNVMix [12]. Paired end libraries were aligned using a modified version of BWA to a reference consisting of the human genome reference (NCBI build 36, hg18) and a database of known exon-exon junctions obtained from different annotation databases (Ensembl [13], RefSeq [14], AceView [15]). Sequences representing exon-exon junctions were designed to require at least a 4 base pair overlap for split-reads. Considering a read length of 50 base pairs, 46 base pairs on either side of the exon-exon junction were concatenated to represent each exon-exon junction. These alignments were used as the information source for differential expression analysis as well as measurement of exon abundance and intron usage. Additionally, de novo splice sites were predicted using HMMSplicer [16]; these predictions served to determine the effects of splice-site mutations in the transcriptome.

Alignment of RNA-Seq short reads was done using a modified version of BWA (base version 0.5.5 [17], patch unpublished). In order to apply the Burrows-Wheeler transform [18], upon which BWA is built, the reference sequence has to be concatenated into one string. When including exon-exon junction sequences to the alignment reference, these are appended to this string in order for the aligner to consider them as part of the target space. This procedure can cause mapping problems when working with paired end tags. In the last stages of alignment, BWA takes advantage of paired information to increase alignment specificity. When reads map to multiple places, whether with zero or more mismatches, pairing information can be used to determine the best alignment position for the pair. In order to optimize this process, BWA measures the distance between the reads of each pair at each mapping position using the concatenated reference as a base and comparing this distance to the estimated fragment size. Alignment positions in which the distance between the reads (too short or too long) or their alignment orientation do not match are discouraged. Given that the exon-exon junction sequences are appended to the genome reference, any pair where one or both reads map to a junction are broken due to the concatenated distance between them. Our modifications make BWA aware of the true position of each exon-exon junction in the reference sequence, thus eliminating some of the mapping issues described above.

Non synonymous SNVs and indels predicted from RNASeq screened against dbSNP and 1000genomes data that were recurrent by gene in at least 3 cases were submitted for validation using resequencing methods described as above for exome capture libraries.

**Gene expression analysis**   Gene expression of Ensembl54 gene models was estimated by the number of fragments that covered them. A fragment was assigned to a gene when at least one read in the read-pair overlapped an exon in the gene, reads with BWA-assigned mapping quality lower than 10 were discarded. In order to make expression values comparable across libraries, we applied quantile normalization [19] of fragment counts in log2 space. Akin to microarray analysis, quantile normalization has been proposed as a viable method for RNA-Seq data that effectively protects against biased estimates due to over-representation of highly expressed genes [20].

**PAM50 classification**   Normalized expression values (as described above) were used classify the triple negative breast cancer samples into two groups (*Basal* and *Other*) using the 50-gene classifier described by Parker et al. [21]. Of the 80 triple negative tumours 48 were classified as *basal*: and 32 as *other*; the latter group included any sample labeled by PAM50 as LuminalA, LuminalB, Her2 or Normal-like. Unsupervised hierarchical clustering (method: Ward, distance: spearman correlation) was performed to assess the reproducibility of PAM50 groupings. At $k = 3$ clusters there is a strong group of 16 samples labeled *Other*, a group formed by 6 libraries obtained from normal tissue and a group of 64 samples with mixed labels; at $k = 4$ this last group is split into a strong *Basal* cluster ($n = 28$) and a mixed one ($n = 36$). Unsupervised hierarchical clustering thus independently supported the primary classification by PAM50 into a relatively homogenous "basal" expression group and a more heterogenous "other" TNBC group.

**Gene fusion analysis and PCR validation**   Gene fusions from paired end RNASeq data were predicted using *deFuse* v0.3.6 [22]. The deFuse predictions were filtered by thresholding on the deFuse probability at a value of 0.5, corresponding to a false positive rate of 0.93 and a false negative rate of 0.15. We identified, from the filtered deFuse predictions, those fusions that were corroborated by a genomic breakpoint in the Affymetrix SNP 6.0 copy number analysis. A fusion was considered corroborated by a genomic breakpoint if that breakpoint was predicted to occur within the boundaries of one of the genes involved in the gene fusion. From the list of genomic breakpoint corroborated fusions, we selected 64 events for validation.

Where available, we sought to identify evidence of genomic rearrangement to corroborate fusion predictions. We performed a targeted search for supporting Whole Genome Shotgun Sequencing (WGSS) read pairs in SOLiD data for 15 tumour genomes. For each fusion, a paired end alignment was considered supporting if it aligned less than 200,000 nucleotides downstream of the predicted fusion boundary in each gene. Supporting paired end alignments for each fusion were clustered, and the largest cluster for each was identified. Sequences for each breakpoint were assembled from paired end alignments as described previously [23]. In total, we identified 29 genomic breakpoints supporting gene fusions. In addition, 5 tumour genomes were sequenced using paired end whole genome shotgun sequencing at low coverage ($\sim$8x) using Illumina GA$_{ii}$. For these genomes, we used comrad v0.0.1 [23] (`http://compbio.cs.sfu.ca`) to identify fusions with supporting genomic breakpoints.

The expression of the fusion transcripts were validated by RT-PCR followed by Sanger sequencing. In brief, reverse transcription was performed using SuperScript III (Invitrogen Canada Inc., Burlington, ON) and 200 ng of RNA from the source tumour. PCR was carried out using Finnzymes PhusionTM High Fidelity DNA polymerase (New England BioLabs Inc., Pickering, Ontario) as well as forward T7 tagged and M13 reverse tagged PCR primers specifically designed to each predicted rearrangement. Products were visualized on 1.5% agarose gel prior to submission for Sanger sequencing at Macrogen (Rockville, MD, USA).

# 6   Network and pathway analysis using Reactome and EnrichmentMap

{ sec:cytoscape }

For network and pathway enrichment analysis for mutations, differential expression and copy number analysis, we used the Reactome FI [24] Cytoscape [25] plugin in Cytoscape v2.8.1. Significantly enriched pathway gene sets were exported and analysed using EnrichmentMap [26] to determine relationships between pathways. For this analysis, gene sets in GMT format were provided by the Reactome FI authors (Wu, personal communication). Patient-pathway clustering was performed by projecting a patient-gene matrix containing occurrences of mutations in a given gene in a given patient to a patient-pathway matrix

containing occurrences of mutations in a given gene belonging to a pathway (or geneset) in a given patient. This patient-pathway matrix was then clustered using hierarchical agglomerative clustering with Euclidean distance and Ward linkage. We acknowledge the redundancy induced by the many-to-many relationship structure of genes in pathway gene sets, however to preserve biological interpretability and to maintain identifiable gene sets in Reactome, genesets were kept unaltered.

# 7  Affymetrix SNP6 analyses

## 7.1  Normalization of intensities

The Affymetrix SNP6.0 arrays for both tumour and normal samples were all independently normalized using the single array method CRMAv2 [27]. We applied the default settings using the following tags: ACC,ra,-XY,BPN,-XY,AVG,A+B,FLN,-XY. For each sample, allelic-crosstalk calibration, probe sequence effects normalization, probe-level summarization, and PCR fragment length normalization was performed. The following annotation files were used (Feb 14, 2008): Chip definition file (GenomeWideSNP_6.Full.cdf), Unit fragment-length (GenomeWideSNP_6,Full,na24,HB20080214.ufl), Unit genome position (GenomeWideSNP_6,Full,na24,HB

## 7.2  Normalization using a reference

Log ratios were computed for both tumours and normals by normalizing each array independently against a reference. Due to the incomplete pairings of tumours and matched normals, we needed to generate a pooled reference from all the available normals.. However, because we are also attempting to identify CNVs in addition to CNAs in the tumours, using a simple pooled reference (median intensity across all normals for each probe) can "subtract" inherently frequent ($> 50\%$) copy number events in the normals from the tumours. For example, using a reference which contains a region that is deleted in more than half the normal samples will result in an artificially amplified signal when computing the ratio for a true neutral signal in the tumour sample. Therefore, we generated a "masked" reference from the normal dataset that was free of inherent copy number polymorphisms, while accounting for systematic biases.

We started with CRMAv2 normalized intensities of the matched normal samples and computed log ratios using a randomly permuted reference. Next, probe-level discrete copy number calls were made using a 6-state hidden Markov model (described below). For each sample, we identify probes that are predicted as gains or losses. Next, we retrieved the original normalized intensities that we started with and masked out the probes we identified in the previous step. Finally, we computed the median (pooling) intensity value for each probe across all samples using only the remaining values after masking. The result is a single vector of median values, which makes up the masked reference. In this approach, the masked reference improves analysis of germline CNVs in the tumours by serving as a true base-line vector of intensities for log ratio calculations and allowing all true alteration signals in the data to be detected.

## 7.3  Segmentation and copy number analysis of 28 Affymetrix SNP6 normal samples using a 6-state HMM

We applied the 6-state HMM to 28 normal samples to detect germline CNVs by returning segment and probe-level copy number. The 6-state HMM is a modified version of CNA-HMMer [28] that that has been adapted to analyzing high-density genotyping arrays such as the Affymetrix SNP6.0 platform and whole genome shotgun data. Another significant extension is the inclusion of additional copy number states, which offers a more intuitive interpretation of DNA dosage in cancer. The new discrete copy number state space included neutral (NEUT), homozygous (HOMD) and hemizygous (HETD) deletions, gain (GAIN), amplification, (AMP) and high-level amplification (HLAMP).

$$K_{CNA} = \{HOMD, HETD, NEUT, GAIN, AMP, HLAMP\} \tag{1}$$ {eq:Kcna}

## 7.4 Segmentation and copy number analysis of 104 Affymetrix SNP6 tumour samples using HMM-Dosage

We performed copy number analysis on the 104 tumour samples using HMM-Dosage (Ha et al, manuscript in preparation). The alrogithm is designed to detect and distinguish the complete set of somatic and germline copy number events in cancer genomes interrogated by SNP array data. This model extends the 6-state HMM (above) by using 5 additional CNV states to represent the analogous copy number status

$$K_{CNV} = \{CNVHOMD, CNVHETD, CNVGAIN, CNVAMP, CNVHLAMP\}. \qquad (2)$$ {eq:Kcnv}

The model performs segmentation on log ratios of intensity data and assigns a discrete copy number status $k_t \in \{K_{CNA}, K_{CNV}\}$ to the each latent variable, $Z_t$, for all probes $t \in \{1 \ldots T\}$. Log ratios, $Y = \{y_1, \ldots, y_T\}$, are observed values modeled using Student's-t densities in the emission component of the HMM, conditional on $Z_t = k_t$. Spatial information is captured using a non-stationary transition matrix, $A_t$, that encodes position-specific probabilities of CNVs and CNAs at each probe $t$. In order to distinguish between somatic copy number alterations (CNAs) and germline copy number variants (CNVs), HMM-Dosage probabilistically incorporates CNV information as a prior to the transition matrix.

The CNV prior was computed as probe-level CNV frequencies by combining (weighted averaging) 2 datasets: 482 normal samples from METABRIC [29] and an external dataset of 450 HapMap normal samples whose CNVs were predicted by [30]. Because the HapMap dataset contains males, we excluded chrX frequencies and simply used the chrX frequencies from the 482 normals.

The parameters of the Student's-t distributions are unobserved and estimated using the expectation maximization (EM) algorithm for each sample, independently. Initial parameters for the $K_{CNA}$ states of the Student's-t mixture were empirically determined using 45 SNP6 breast cancer cell line samples in the COSMIC [31] dataset. First, independently for each cell-line sample, we fit the log ratios to a 6-state Gaussian mixture model using EM. Subsequently, the converged Gaussian mixture parameters are used as initial parameters for fitting a 6-state Student's-t mixtures in the 6-state HMM (again, using EM). Finally, we averaged the converged Student's-t parameters for each state across the 45 cell lines. These averaged parameters became the initial $K_{CNA}$ parameters to Student's-t emission of HMM-Dosage. Initial parameters for $K_{CNV}$ were set by hand to theoretical values $log2([0.5, 1, 3, 5, 7]/2)$.

The optimal state sequence, $Z_{1:T}$, which represents the copy number prediction for each probe, is computed using the Viterbi algorithm.

In general, HMM-Dosage was tuned to be conservative in calling CNVs as to avoid misclassifying the more important CNA calls. Therefore, the results were post-processed by comparing against the CNV map (derived from 482 normals and HapMap270) such that CNA segments that had at least 25% reciprocal genomic overlap with a segment of the same copy number state in the CNV map were converted to germline state.

## 7.5 Genomic instability

For each of the 80 genomes analyzed by Affymetrix SNP6 arrays that have RNA-seq data, genomic instability (GI) was computed as the proportion of genome altered by basepair. For this analysis, we calculated the total length of all CNA events (HOMD, HETD, GAIN, AMP, HLAMP) by using the segmentation boundaries predicted by HMM-Dosage. This value is then divided by the total length of all predicted segments, including NEUT. This calculation is also done for each CNA type to reveal the relative contribution of each alteration event type to GI.

We also compared GI between basal and non-basal PAM50 subtypes. The basal subtype shows higher GI for all CNA types.

## 7.6 Gene alterations

In order to identify the genes that are altered by copy number changes, we searched for overlap of segments with gene regions. The gene annotations and coordinates are given by Ensembl 54 (hg18). The genes selected for analysis are those annotated as protein-coding, totaling to 20878 genes.

We generate two types of patient-by-gene copy number matrices to capture two perspectives of gene alterations. The first type is the 'call' matrix, $C \in \mathbb{Z}^{P \times G}$, which is populated with values representing discrete copy number state calls. The second is the log ratio matrix, which contains segment median log ratios. For this latter matrix, genes overlapping multiple segments are summarized using two different approaches to produce $L_a \in \mathbb{R}^{P \times G}$ and $L_s \in \mathbb{R}^{P \times G}$. For each patient $p \in P$ and each gene $g \in G$, we identify segment $s$ that overlaps $g$ and assign $C(p, g)$ with the copy number state of $s$ and $L_a(p, g)$ and $L_s(p, g)$ with the median log ratio of the probes in $s$.

If gene $g$ overlaps or is broken by a set of segments, $S = s_1, \ldots, s_k$, where $k \geq 2$, we assign $C(p, g)$ with the copy number state of the segment having max *severity* based on the binary relation (Equation 5 and 6). For $L_a(p, g)$, the weighted sum of the median logRs (Equation 3) is used while for $L_s(p, g)$, the logR corresponding to the most severe segment is used (Equation 4).

$$L_a(p, g) = \left( \sum_{s \in S} length(s) * medianLogR(s) \right) / \sum_{s \in S} length(s) \tag{3}$$

$$L_s(p, g) = logRvalueOf \left( \underset{s \in S}{argmax}\{severity(CNstateOf(s))\} \right) \tag{4}$$

$$C(p, g) = CNstateOf \left( \underset{s \in S}{argmax}\{severity(CNstateOf(s))\} \right) \tag{5}$$

$$
\begin{aligned}
severity = \{ &(`NEUT', 0), \\
&(`HOMD', 8), (`HETD', 6), \\
&(`GAIN', 5), (`AMP', 6), (`HLAMP', 8), \\
&(`CNVHOMD', 4), (`CNVHETD', 3), \\
&(`CNVGAIN', 2), (`CNVAMP', 3), (`CNVHLAMP', 4)\}
\end{aligned} \tag{6}
$$

## 7.7 Subtype-specific analysis of copy number

A $\mathcal{X}^2$ (Chi-Square) test of independence was used to determine subtype association with copy number. For each gene $g$, copy number of the patients were divided into 3 levels: deletion (HOMD & HETD), gain (GAIN, AMP, HLAMP) and neutral. For the omnibus test, patients were divided into 2 levels (subtypes), Basal and non-Basal. No significant genes were found after using Bonferroni adjusted p-value cutoff of 0.1.

## 7.8 Outlier analysis

To identify genes whose expression was driven in cis by extreme copy number events, we fit Gaussian distributions to the 80 expression values of a gene across the patients in the cohort (for which RNASeq data was available), using the maximum likelihood estimates of the mean and variance. For genes with multiple probes, we summarized the expression level using the median. This resulted in 20878 Gaussian distributions, each corresponding to the expression profile of a gene.

We then sought to capture gene-patient $(g, p)$ events where large amplitude copy number events (such as putative homozygous deletions and high level amplifications) drives expression into the extreme tails of the expression profile. We selected events in 5% left-tail or 5% right-tail of the expression distributions. Events in the right tail whose copy number state was either amplification or high-level amplification (HLAMP) by HMM-Dosage and those events in the left tail whose copy number state was homozygous deletion (HOMD) were reported.

# 8   WGSS analysis for copy number

The following analysis is implemented into an R package called HMMcopy available from `http://compbio.bccrc.ca/software/hmmcopy/`.

## 8.1 Segmentation and copy number analysis of 15 WGSS tumour samples using 6-state HMM

The WGSS-derived copy number results used as input in APOLLOH (see below) were generated using a modified version of a paired tumour-normal strategy described previously in [32] and [9]. The genome was divided into fixed windows of 1kb, and read depth is extract for each window in the tumour and normal. In this study, we applied two additional preprocessing steps prior to segmentation in order to achieve more accurate copy number estimates. First, we applied a filter to remove repetitive regions that are highly mappable. Second, we corrected GC content bias to remove wave-like patterns in the tumour and normal, separately, using a loess curve fit between GC content and read depth. A log ratio is computed for each window by computing proportion between the GC corrected tumour value and GC corrected normal value.

### Divide the genome into fixed genomic windows

First, the genome was divided into fixed genomic windows of 1kb, reducing the analysis to a set of approximately 2.5 to 3 million loci, $\mathbf{R}$.

### Extract read depth for normal and tumour

Next, separately for the normal and tumour genomes of each patient, we extracted the total read depth for each window in $\mathbf{R}$ using BAMtools [33], resulting in a vector of read count data from the normal, $\mathbf{N_R} = (n_1, \ldots, n_{|R|})$ and tumour, $\mathbf{T_R} = (t_1, \ldots, t_{|R|})$, respectively.

### Removing windows that are highly mappable

Using "ENCODE Duke Uniqueness of 35bp sequences" track from UCSC (`http://genome.ucsc.edu/cgi-bin/hgTrackUi?db=hg18&g=wgEncodeMapability`), we filtered windows that were within repetitive regions, resulting in being highly mappable by aligners. Windows that had mappability score of $\geq 0.9$ were excluded. This removed extreme amplified positions which would have otherwise posed as confounding outliers in downstream segmentation analysis.

### GC content correction of normal and tumour read counts

We performed GC content bias correction to the tumour and normal of each patient separately. We applied a global loess fit between GC content and read depth for windows in $\mathbf{R}$. Due to computational restrictions of fitting 3 million data points, we further excluded outlier windows based on read depths being in the upper and lower 1% quantile, and randomly sampled 20,000 of the remaining windows for generating the loess curve. Finally, the read depth of all windows, $\mathbf{N_R}$ and $\mathbf{T_R}$, are corrected by scaling the observed value by the loess fitted value (Equation 7).

$$corrected\ read\ depth = \frac{observed\ read\ depth}{loess\ fitted\ value} \tag{7}$$

{eq:loess`correct}

### Normalizing copy number in tumours

The GC corrected normal $\mathbf{N_R}$ and tumour $\mathbf{T_R}$ counts are normalized independently to generate $\bar{\mathbf{N}}_\mathbf{R} = (\bar{n}_1, \ldots, \bar{n}_{|R|})$ and $\bar{\mathbf{T}}_\mathbf{R} = (\bar{t}_1, \ldots, \bar{t}_{|R|})$, respectively, where $\bar{n}_i = \frac{n_i}{\sum_j n_j}$ and $\bar{t}_i = \frac{t_i}{\sum_j t_j}$, $i \in \{1, \ldots, |R|\}$. To obtain the final tumour copy number observed at each loci $r \in R$, we applied another normalization step by taking the log2 ratio between tumour and normal copy number, $\mathbf{T_N(i)} = log_2\left(\frac{\bar{t}_i}{\bar{n}_i}\right)$, $i \in \{1, \ldots, |R|\}$.

**Segmentation and copy number prediction via HMM**

The 6-state version of HMM-Dosage was used to segment the input data $\mathbf{T_N}$. Initialization and hyperparameters for the prior means of the Student's-t distribution used for this analysis were $log2(([1, 1.4, 2, 2.7, 3, 4.5])/2)$. Because the normalized the tumour counts with the match normal read count naturally excludes germline CNVs, using the 11-state HMM-Dosage to distinguish CNA from CNV was not required.

## 8.2 Segment-based correlation of WGSS and SNP6 predicted copy number

Copy number results from WGSS data was used to help validate HOMD and HLAMP predictions in the Affymetrix SNP6 data. In the 15 cases, for each HOMD and HLAMP predicted segment, $s$, boundaries $s_{start}$ and $s_{end}$ and segment median $a_s$ are extracted from the SNP6 results. We compute the median log ratio copy number $w_s$ across the set of regions $i \in R$ that overlap these boundaries $s_{start}$ and $s_{end}$.

$$w_s = median(T_N(r)), \ \{r|r \in R, r_{start} \geq s_{start}, r_{end} \leq s_{end}\}$$

Spearman correlation is computed on $w_{1:S}$ and $a_{1:S}$ where $S$ is the total number HOMD/HLAMP segments in the SNP6 results.

# 9 Analysis of LOH

We determined loss of heterozygosity (LOH) in the 15 WGSS and 50 exome samples using a probabilistic model packaged in the software called APOLLOH (Ha, et al manuscript in preparation). Preprocessing of the data is performed in a paired tumour-normal approach by interrogating informative heterozygous loci identified in the normal genomes using UnifiedGenotyper from GATK [34]. Next, the tumour genome allelic counts in the same loci are input into APOLLOH which performs segmentation of the allelic ratios into regions of LOH, allele-specific copy number amplification (ASCNA) and heterozygosity (HET). Copy number was also input into APOLLOH to enhance predictions within regions that deviate from copy neutral. For the WGSS samples, we used copy number results predicted by the previously described depth-based WGSS approach. For the exome samples, we used copy number results from the SNP6 genotyping arrays. Software for APOLLOH is available from the authors' website at : `http://compbio.bccrc.ca/software/apolloh/` .

# 10 Human regulatory region analysis

Analysis of 323 high-confidence somatic mutations in non-coding space showed that 72 were found to fall within experimentally determined human regulatory regions as described in [35]. For each variation contained within a ChIP-seq-defined region, the matrix for the corresponding TF was used to predict putative TFBSs overlapping the variation (reviewed in [36]). Two TFBS scores were generated, one corresponding to the reference genome sequence and the other corresponding to the observed cancer variation. A relative matrix score threshold of 80% was applied to define putative TFBSs. As the variations were selected to overlap regions of known transcription factor binding, we generally expect the reference allele to score higher for functionally significant changes (i.e. to go from a functional TFBS to a non-functional one). We called regulatory variations exceeding the 80% threshold as putatively impactful to transcription factor binding. In the TN breast cancer set, 2 such somatic mutations were predicted to affect MITF binding and 6 somatic mutations were predicted to affect RB binding.

# 11 DriverNet: Integrated analysis of copy number, mutations and gene expression data

To nominate mutations predicted to have significant impact on expression networks, we used the 681 Reactome FI gene sets containing 7521 genes to construct an 'influence graph' where genesets were used

to form 'cliques' in the graph. We then constructed a bipartite graph such that each gene that was in the mutation list from sequencing or HOMD, HLAMP list from CNA analysis was placed in one partition of the graph and the outlying expression genes (as described above) was placed in the other partition. Edges between nodes in the two partitions were drawn according to the following criteria:

- a mutation and an outlying expression event co-occurred in the same tumour

- an edge between the two genes was present in the influence graph

With the influence graph constructed, we then sought to rank the genes in the 'mutation' partition such that nodes explaining the highest number of outlying expression expression events (by greedy algorithm) were nominated as putative driver mutations. The greedy algorithm is an efficient approximation to the combinatorial optimization problem induced by the bipartite graph construction that we wish to solve in order to predict the fewest number of mutations that disrupt the maximum number of expression signatures.

The statistical significance of the candidate genes are assessed using a randomization framework. Given the original dataset $S$, the algorithm is run on randomly generated datasets N times $(S_1, \ldots, S_N)$ and then the results on real data are assessed to see if they are significantly different from the randomized datasets. This in an indirect way of perturbing the bipartite graph corresponding to the original problem. To generate the random datasets, we keep the contents of patient-mutation, $M$, and patient-outlier, $G'$, matrices the same but replace the gene symbols with a randomly selected set of genes from the Ensmbl 54 protein-coding gene list. Using the same influence graph, the algorithm is run on the new patient-mutation, $M_1...M_M$, and patient-outlier, $G'_1...G'_M$, matrices.

Suppose $D$ is the result of driver mutation discovery algorithm on the input $U$. $D$ contains a ranked list of driver genes with their corresponding node coverage in the bipartite graph, $\mathcal{B}$. The statistical significance of a gene $g \in D$ with a corresponding node coverage, $COV_g$, is the fraction of times that we observe driver genes in our random data runs $D_i$, with node coverage more than $COV_g$:

$$\text{pvalue}(g) = \frac{\sum_{i=1}^{N} \sum_{j=1}^{M} \delta[COV_{gj} > COV_g]}{N * \sum_{i=1}^{N} M_i} \tag{8}$$

where $M_i$ is the number of drivers identified in run $n$ of the algorithm, $N$ is the total number of times that the algorithm on random datasets were run.

Software for this analysis called *DriverNet* (Haffari et al., manuscript in preparation) is available from the authors' website: `http://compbio.bccrc.ca/software/drivernet`. Pathway analysis of the top 50 genes output by the greedy algorithm was performed using Reactome as described above.

# 12    Clonal frequency analysis

{ sec:clonal }

We used a hierachical Bayesian model which accounts for uncertainty about genotype to estimate the frequency of cells with a given mutation. We imagine each mutation divides the population of cells into two classes; those without the mutation we call the reference population and those with we refer to as the variant population. We treat the reference population as though it has fixed genotype of AA. We model the genotype of the variant population at a given locus using one of 27 states allowing for copy numbers from 1-6 and differing zygosity status. We use information about copy number and LOH status derived from APOLLOH exome analysis to inform or prior beliefs on the genotype of the variant population. To allow multiple mutations to have the same frequency, we use a Dirichlet process prior on the class frequencies. Sites which are clustered by this model have cellular frequencies which are sufficiently similar that they are better explained as originating from a single class of cells. Inference for this model is performed using Markov chain Monte Carlo sampler. A full description of the hierarchical Bayesian model shown as a probabilistic graphical model is shown in Figure S10. Input to the model is a set of mutations and their allele abundance (reference, variant). In addition, the copy number and LOH status at the position of each mutation is given. The output of the model is a mutation co-occurrence matrix containing the proportion of MCMC samples in which each pair of mutations were grouped together, and the clonal frequency estimation for each of the MCMC samples, thus allowing inference of a posterior distribution of clonal frequency estimates.

Software implementing this approach called *PyClone* (Roth et al., manuscript in preparation) is available upon request from the authors. Cases with 10 more more mutations were included in the clonal frequency analysis. This led to estimates for 2109 mutations distributed across 54 cases. We clustered the mutation co-occurrence matrix for each case using hierarchical clustering and deterministically assessed the groupings of mutations into 'clonal clusters' using the R library dynamicTreeCut [37]. Clonal frequency distributions for each mutation were then plotted as in Figure 4 according to these groupings.

Each pathway in Reactome was also assessed for skewing of clonal frequency by Wilcoxon test. Clonal frequency estimates for each case were first adjusted by normalizing by the maximum clonal frequency value to remove the effect of normal cell contamination. The adjusted clonal frequencies were then used to assess whether the distributions of the mutations falling into a given pathway were different than the background distribution of all clonal frequency estimates using the R implementation of the Wilcoxon test `wilcox.test`. Resultant p-values were then adjusted using `p.adjust` with `method="BH"`.

# References

[1] Wiegand, K. C. *et al.* Arid1a mutations in endometriosis-associated ovarian carcinomas. *N Engl J Med* **363**, 1532–1543 (2010).

[2] Gnirke, A. *et al.* Solution hybrid selection with ultra-long oligonucleotides for massively parallel targeted sequencing. *Nat Biotechnol* **27**, 182–189 (2009).

[3] Morin, R. D. *et al.* Frequent mutation of histone-modifying genes in non-hodgkin lymphoma. *Nature* **476**, 298–303 (2011).

[4] McKernan, K. J. *et al.* Sequence and structural variation in a human genome uncovered by short-read, massively parallel ligation sequencing using two-base encoding. *Genome Research* **19**, 1527–1541 (2009).

[5] Li, H., Ruan, J. & Durbin, R. Mapping short dna sequencing reads and calling variants using mapping quality scores. *Genome Res* **18**, 1851–1858 (2008).

[6] Roth, A. *et al.* Jointsnvmix : A probabilistic model for accurate detection of somatic mutations in normal/tumour paired next generation sequencing data. *Bioinformatics* (2012).

[7] Ding, J. *et al.* Feature-based classifiers for somatic mutation detection in tumour-normal paired sequencing data. *Bioinformatics* **28**, 167–175 (2012).

[8] Reva, B., Antipin, Y. & Sander, C. Predicting the Functional Impact of Protein Mutations: Application to Cancer Genomics. *Nucleic Acids Research* (2011, in press).

[9] Shah, S. P. *et al.* Mutational evolution in a lobular breast tumour profiled at single nucleotide resolution. *Nature* **461**, 809–13 (2009).

[10] Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078–2079 (2009).

[11] Wu, T. D. & Nacu, S. Fast and snp-tolerant detection of complex variants and splicing in short reads. *Bioinformatics* **26**, 873–881 (2010).

[12] Goya, R. *et al.* SNVMix: predicting single nucleotide variants from next-generation sequencing of tumors. *Bioinformatics* **26**, 730–736 (2010).

[13] Flicek, P. *et al.* Ensembl's 10th year. *Nucl. Acids Res.* **38**, D557–562 (2010).

[14] Pruitt, K. D., Tatusova, T. & Maglott, D. R. NCBI reference sequences (RefSeq): a curated non-redundant sequence database of genomes, transcripts and proteins. *Nucl. Acids Res.* **35**, D61–65 (2007).

[15] Thierry-Mieg, D. & Thierry-Mieg, J. AceView: a comprehensive cDNA-supported gene and transcripts annotation. *Genome Biology* **7**, S12 (2006).

[16] Dimon, M. T., Sorber, K. & DeRisi, J. L. HMMSplicer: a tool for efficient and sensitive discovery of known and novel splice junctions in RNA-Seq data. *PLoS ONE* **5**, e13875 (2010).

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

[18] Burrows, M. & Wheeler, D. J. A block-sorting lossless data compression algorithm (1994).

[19] Bolstad, B., Irizarry, R., Astrand, M. & Speed, T. A comparison of normalization methods for high density oligonucleotide array data based on variance and bias. *Bioinformatics* **19**, 185 –193 (2003).

[20] Bullard, J., Purdom, E., Hansen, K. & Dudoit, S. Evaluation of statistical methods for normalization and differential expression in mRNA-Seq experiments. *BMC Bioinformatics* **11**, 94 (2010).

[21] Parker, J. S. *et al.* Supervised risk predictor of breast cancer based on intrinsic subtypes. *Journal of Clinical Oncology* **27**, 1160 –1167 (2009).

[22] McPherson, A. *et al.* deFuse: An Algorithm for Gene Fusion Discovery in Tumor RNA-Seq Data. *PLoS Comput Biol* **7** (2011).

[23] McPherson, A. *et al.* Comrad: detection of expressed rearrangements by integrated analysis of rna-seq and low coverage genome sequence data. *Bioinformatics* **27**, 1481–1488 (2011).

[24] Wu, G., Feng, X. & Stein, L. A human functional protein interaction network and its application to cancer data analysis. *Genome Biology* **11**, R53 (2010).

[25] Smoot, M. E., Ono, K., Ruscheinski, J., Wang, P. L. & Ideker, T. Cytoscape 2.8: new features for data integration and network visualization. *Bioinformatics* **27**, 431–432 (2011).

[26] Merico, D., Isserlin, R., Stueker, O., Emili, A. & Bader, G. D. Enrichment map: a network-based method for gene-set enrichment visualization and interpretation. *PLoS ONE* **5**, e13984 (2010).

[27] Bengtsson, H., Wirapati, P. & Speed, T. P. A single-array preprocessing method for estimating full-resolution raw copy numbers from all affymetrix genotyping arrays including genomewidesnp 5 & 6. *Bioinformatics* **25**, 2149–56 (2009).

[28] Shah, S. P. *et al.* Integrating copy number polymorphisms into array cgh analysis using a robust hmm. *Bioinformatics* **22**, e431–9 (2006).

[29] Curtis, C. *et al.* The genomic and transcriptomic architecture of 2,000 breast tumours reveals novel subgroups. *Nature* (2012, accepted for publication).

[30] Conrad, D. F. *et al.* Origins and functional impact of copy number variation in the human genome. *Nature* **464**, 704–12 (2010).

[31] Bignell, G. R. *et al.* Signatures of mutation and selection in the cancer genome. *Nature* **463**, 893–898 (2010).

[32] Chiang, D. Y. *et al.* High-resolution mapping of copy-number alterations with massively parallel sequencing. *Nat Methods* **6**, 99–103 (2009).

[33] Barnett, D. W., Garrison, E. K., Quinlan, A. R., Strömberg, M. P. & Marth, G. T. Bamtools: a c++ api and toolkit for analyzing and managing bam files. *Bioinformatics* **27**, 1691–1692 (2011).

[34] McKenna, A. *et al.* The genome analysis toolkit: a mapreduce framework for analyzing next-generation dna sequencing data. *Genome Res* **20**, 1297–1303 (2010).

[35] Chicas, A. *et al.* Dissecting the unique role of the retinoblastoma tumor suppressor during cellular senescence. *Cancer Cell* **17**, 376–387 (2010).

[36] Worsley-Hunt, R., Bernard, V. & Wasserman, W. W. Identification of cis-regulatory sequence variations in individual genome sequences. *Genome Med* **3**, 65–65 (2011).

[37] Langfelder, P., Zhang, B. & Horvath, S. "defining clusters from a hierarchical cluster tree: the dynamic tree cut package for r. *Bioinformatics* **24**, 719–720 (2007).

# Supplemental Tables for Shah et al: The clonal and mutational evolution spectrum of primary triple negative breast cancers

This file contains the legends describing Tables S1-17. All tables are available as an archived zip file in the Supplemental Data from the Nature website.

Table 1: Summary of data generation for 54 tumour and normal DNA exome libraries and 80 RNASeq transcriptome libraries. RNAseq and exome libraries were 50bp paired end Illumina sequence data. For RNASeq data, the number of aligned reads in exons, introns and intergenic regions of the genome are listed with three mapping quality categories. Mean, median and standard deviation of targeted exon coverage is listed for each exome library. Total acquired data for 15 SOLiD whole genome shotgun tumour/normal pairs are also listed.

Table 2: Summary of clinical characteristics of the patients. A full data dictionary describing all clinical variables is presented (Sheet 2)

Table 3: Omnibus table of validated somatic SNVs. All mutations were revalidated with deep amplicon resequencing. There were 2414 somatic point mutations that revalidated. Mutations were annotated with MutationAssessor (http://mutationassessor.org/, [1]). Column descriptions are as follows: **id**: unique patient/genome coordinate (NCBI build36, hg18); **Patient.ID**: case id; **Chr: chromosome**; **Coord**: chromosomal position; **Mutation**: nucleotide substitution; **RG.var.type**: missense, nonsense or silent mutation; **Func..Impact**: mutationassessor Functional impact class; **FI.score**: mutationassessor Functional impact score; **AA.variant**: amino acid substitution; **Gene**: gene name (HGNC); **Location**: chromosomal band; **Uniprot**: uniprot id; **Refseq**: refseq id; **mutations.in.COSMIC.position**: previously reported mutations at the same genomic coordinate in the COSMIC database; **types.of.cancer.in.COSMIC.position**: description of previously reported COSMIC mutations; **gene.s.known.role.in.cancer**: cancer gene census; **regions.position**: affected protein/domain localisation; **Ref.count.Tumour.**: resquencing reads from the tumour DNA matching the reference; **Alt.count.Tumour.**: resquencing reads from the tumour DNA matching the variant; **Total.count..Tumour.**: total resequencing reads from the tumour DNA; **Alt.frequency..Tumour.**: proportion of resequencing reads that are variant; **Binomial.test.p.value..Tumour.**: Binomial exact test p-value; **Binomial.test.corrected.p.value..Tumour.**: Binomial exact test corrected p-value; **Binomial.test.status..Tumour.**: Binomial exact test corrected p-value < 0.01?; **Ref.count..Normal.**: resquencing reads from the matched normal DNA matching the reference; **Alt.count..Normal.Total.count..Normal.**: resquencing reads from the normal DNA matching the variant; **Alt.frequency..Normal.**: proportion of resequencing reads matching the variant; **Binomial.test.p.value..Normal.**: Binomial exact test p-value; **Binomial.test.corrected.p.value.Normal.**: Binomial exact test corrected p-value; **Binomial.test.status.Normal.**: Binomial exact test corrected p-value < 0.01?; **SVNMix2.Call.RNASeq**: most probable genotype from SNVMix2 [2] (0 - not expressed, 1 - wildtype expressed, 2 - heterozygously expressed, 3 - homozygously expressed) ; **Ref.RNASeq.**: number of reference reads at position from the matching RNASeq data ; **Nonref.RNASeq.**: number of variant reads at position from the matching RNASeq data; **Ens54Gene:** Ensembl v54 gene name ; **Strand**: gene strand ; **CNCall**: copy number state ; **CNLogR**: SNP6.0 median log ratio for the overlapping segment at this position ; **PAM50**: PAM50 class (Basal/Other) for this case ; **Splice.site**: position is a dinucleotide donor/acceptor splice site (TRUE/FALSE); **HMMDosage_Call**: copy number state; **HMMDosage_LogR**: raw continuous copy number value; **OncoSNP_Call**: predicted OncoSNP [3] state ; **OncoSNP_CN**: predicted OncoSNP copy number; **apolloh.loh.ratio**: raw LOH value from APOLLOH analysis (see supplemental methods); **apolloh.call**: LOH call from APOLLOH from exomes; **APOLLOH_WGSS_Call**: LOH call from APOLLOH from genomes;**APOLLOH_WGSS_CN**: copy number used for APOLLOH analysis; **HMMcopy_WGSS_Call**: copy number inferred from genome data; **HMMcopy_WGSS_CN**: normalized logR input for genome-based copy number analysis; **mut_freq_GBM**: frequency of mutation in gene in TCGA GBM data; **mut_freq_HGS**: frequency of mutation in gene in TCGA HGS data; **clonality**: estimation of clonal frequency.

Table 4: Validated SNVs at dinucleotide donor/accptor splice sites (based on Ensembl v54 annotations) inducing alternative splicing or intron retention in the RNASeq data: legends as per Table S3.

Table 5: Validated somatic insertions and deletions. Scores are as predicted by `samtools pileup` [4].

Table 6: Table of CNA segments for 104 Affymetrix SNP6.0 arrays, predicted as described in Supplemental methods. Columns are CaseID; Chromosome; segment start coordinate (NCBI build 36); segment end coordinate; number of probes in segment; median logratio of segment; HMMDosage state (somatic states are HOMD: homozygous deletion, HETD: hemizygous deletion: GAIN: single copy gain, AMP: amplification, HLAMP: high-level amplification (see Supplemental methods); HGNC gene names of genes contained in the segment. The first six columns can be extracted and viewed in the Integrated Genomics Viewer Software (http://www.broadinstitute.org/igv/).

Table 7: Summary of 517 CNA breakpoint-corroborated gene fusions predicted by deFuse [5]. In addition to the columns described at http://compbio.bccrc.ca, 2 columns (gene_1_broken, gene_2_broken) identify which genes are interrupted by a CNV breakpoint, and 2 columns (transition_type_1, transition_type_2) identify the type of the transition.

Table 8: Summary of genomic breakpoints from a targeted search for supporting paired end reads in 15 SOLiD paired end tumour genomes. Shown are the case identifier, genes involved in the fusion, breakpoint sequence, and number of supporting paired end reads.

{ stab:solid˙genome˙breaks }

Table 9: Integrated genome/transcriptome analysis for expressed gene fusions with corroborated genome rearrangements from 5 tumours (SA029, SA030, SA031, SA051, SA054) for which both RNASeq and low-pass whole genome shotgun data generated with Illumina paired end sequencing was obtained. Analysis was performed using Comrad (http://fusioncomrad.sf.net). A description of the table format can be found at http://fusioncomrad.sf.net

{ stab:comrad˙results }

Table 10: A subset of 35/63 selected gene fusion candidates from Table S7 that were validated with PCR.

{ stab:validated˙fusions }

Table 11: 72 non-coding mutations falling in regulatory regions showing 23 mutations in Rb transcription factor binding sites and relative scores of predicted Rb binding in wildtype vs mutant tfbs sequences

{ stab:regregions }

Table 12: Mutation frequency by gene for somatic SNV, indel and splice site mutations

{ stab:snv˙indel˙gene˙freque

Table 13: Mutation significance analysis showing genes predicted to be under selection through analysis of background synonymous mutations and non-synonymous mutations. All genes were analysed using methodology outlined in Youn and Simon [6].

{ stab:mutsig }

Table 14: Enriched pathways from the network computed from recurrent somatically mutated genes in the Reactome functional interaction map. Enriched pathways with FDR <= 0.001 are included.

Table 15: Analysis of the top somatically aberrated genes (by node degree) connected (by Reactome gene sets) to genes that exhibited outlying expression from their population level distributions as computed by driverNet. Columns: Rank: by driverNet algorithm (Supplemental methods); Gene: somatically aberrated gene; gband: chromosomal band containing gene; SNV/Indel: number of cases harbouring an SNV or indel in the gene; HLAMP: number of cases harbouring a predicted high level amplification; HOMD: number of cases harbouring a predicted homozygous deletion; events: number of gene expression outliers (see Supplemental methods) coincident with a genomic aberration and where the outlying gene is connected to the aberrated gene; p-value: statistical significance based on a randomly generated background distribution (Supplemental methods).

Table 16: Comparison of significantly enriched pathways (FDR < 0.005) by Reactome analysis from the top 50 genes predicted by driverNet (Supplemental methods) for TNBC and ovarian high grade serous (HGS) data from the TCGA [7]. The pathway, the initial p-value, the adjusted p-value the driverNet genes and the overlap status are shown, where status=1 indicates TNBC specific pathway, status=2 indicates HGS and TNBC enriched pathways and status=3 indicates HGS specific pathways.

Table 17: Analysis of the distributions of clonal frequency estimates of genes in Reactome pathways. For each set of genes belonging to a pathway, a Wilcoxon test was performed to assess whether the distribution of the clonal frequencies of genes in a pathway was different than the overall background distribution. For each pathway, the name of the Reactome pathway, the geneset of the pathway, the mutated genes that fall into the pathway, the number of such mutation events, the median clonal frequency estimate, the Wilcoxon p-value, and adjusted p-value (Benjamini-Hochberg) are shown.

# References

[1] Reva, B., Antipin, Y. & Sander, C. Predicting the functional impact of protein mutations: application to cancer genomics. *Nucleic Acids Res* **39** (2011).

[2] Goya, R. *et al.* SNVMix: predicting single nucleotide variants from next-generation sequencing of tumors. *Bioinformatics* **26**, 730–736 (2010).

[3] Yau, C. *et al.* A statistical approach for detecting genomic aberrations in heterogeneous tumor samples from single nucleotide polymorphism genotyping data. *Genome Biol* **11** (2010).

[4] Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078–2079 (2009).

[5] McPherson, A. *et al.* deFuse: An Algorithm for Gene Fusion Discovery in Tumor RNA-Seq Data. *PLoS Comput Biol* **7** (2011).

[6] Youn, A. & Simon, R. Identifying cancer driver genes in tumor genome sequencing studies. *Bioinformatics* **27**, 175–181 (2011).

[7] Cancer Genome Atlas Research Network. Integrated genomic analyses of ovarian carcinoma. *Nature* **474**, 609–615 (2011).