**Supplementary Materials**

Online methods
Figures S1, S2, S3, S4, S5, S6, S7, S8, S9, S10, S11
Tables S1, S2, S3, S4, S5, S6, S7, S8, S9, S10, S11, S12
References (supplementary)


**Online methods**

**SEER data**

Surveillance, Epidemiology, and End Results (SEER) data were downloaded from

http://seer.cancer.gov/data/.


**Tumor / normal exon sequencing, DNA copy number, and RNA-sequencing data**

Genomic data was obtained from the TCGA data portal (https://tcga-data.nci.nih.gov/tcga/) and

from additional Broad Institute datasets[39,50]. Informed consent was obtained from all subjects by

local institutional review boards and consents were reviewed by the sequencing centers, per

TCGA guidelines. To visualize gene expression by mutation status for males and females, TCGA

data were downloaded and visualized using cBioPortal (www.cbioportal.org)[51,52].


**Mutation classification**

Tumor-associated DNA variants were called "truncating" if they resulted in nonsense, insertion,

deletion, translation start site, or nonstop mutations. Truncating and likely functional missense

(LOF) included all truncating mutations as well as missense mutations that occurred at the same

site in at least 3 patients (n=419 across the pan-cancer dataset, of which 9 were male-

predominant) or had a Cons 46-vertebrates score (http://ucscbrowser.genap.ca/cgi-

bin/hgTrackUi?db=hg19&g=cons46way)[53] higher than 0.9 indicating high conservation across

species (n=3083, of which 124 were male-predominant).


**Permutation analysis**

601   For each set of events, the probability of observing the number of events seen in males was

602   determined by a series of permutations. Within each tumor type, the probability that any given

603   gene alteration would occur in a male was assumed to be directly related to the number of

604   mutations across all males, divided by the total number of mutations across all samples. This

605   controlled for the number of males in a set as well as the relative mutation rate on chrX in males

606   vs females in the set. For instance, if 60% of coding mutations on chrX in the dataset for a tumor

607   type were in males, then for any given chrX gene, assuming there is no selective advantage for

608   males vs. females, we should see 60% of the events occurring in males. For each gene, the total

609   number of events was counted and a random set of permutations was constructed, in which each

610   event was randomly assigned a 'male' or 'female' value based on this calculated probability. For

611   the pan-cancer analyses, the probability of an event being in a male was calculated on a per

612   tumor type basis, and events were counted and permuted within each tumor type. Once the

613   number of males in the permutation was calculated, it was compared against the observed data

614   to see if the number of males randomly generated in this manner was greater than or equal to

615   that observed in the data. For each gene, 1 million permutations were done, and the fraction of

616   permutations with a greater than or equal number of male events was reported as the P value.

617   Multiple hypothesis correction was then performed on the complete set of genes that were

618   affected by the set of events examined, and genes with FDR <0.1 (by Benjamini-Hochberg false

619   discovery rate correction) were reported as significant. For permutations on truncating and LOF

620   mutations, this probability was calculated based on the coding mutation count of chrX. For

621   permutations on copy-loss events, the frequency of copy-loss events on chrX was used. For the

622   combined LOF/CN loss permutations, the sum of these two was used.

623

624   **Log likelihood ratio test**

625   Assuming that the only factor differentiating the male mutation rate of a gene and the female

626   mutation rate of a gene is the difference in the background mutation rate on chrX (i.e. assuming

24

627 there is no selective bias in males for mutating a gene), then the probability of a male having a

628 mutation in a gene should be directly related to the probability of a female having a mutation in a

629 gene, corrected for the number of copies of chrX in females (n=2) and males (n=1):

630 $$p_{female} = r * p_{male}$$

631 Where r is the F/M ratio of coding mutations across chrX, i.e. if females in a set have twice the

632 number of coding mutations on chrX as males, it can be expected that any given gene on X

633 should be twice as likely to be mutated in a female patient vs. a male patient. For each test, the

634 actual F/M ratio of coding mutations on chrX in the analysis set was calculated and used.

635 Because there is a direct relationship between the probability of a male mutation in this model

636 and a female mutation, we can express the likelihood that the observed data is consistent with

637 this model using a single value of p:

638 $$L_0 = max((p_{male})^m (1-p_{male})^{M-m} (r*p_{male})^f (1-r*p_{male})^{F-f})$$

639 Where $M$ and $F$ are the total number of male and female patients, respectively, and $m$ and $f$ are

640 the number of mutated males and females, respectively. The alternative hypothesis is that there

641 are independent factors affecting males and females, in which case two values of p are needed to

642 calculate the likelihood that the data fits the model:

643 $$L_1 = max((p_{male})^m (1-p_{male})^{M-m} (p_{female})^f (1-p_{female})^{F-f})$$

644 Which can be maximized by using the observed mutation counts:

645 $$L_1 = (m/M)^m (M-m/M)^{M-m} (f/F)^f (F-f/F)^{F-f}$$

646 The log-likelihood ratio (*LLR*) is calculated simply by taking the log of the ratio of these two

647 numbers:

648 $$LLR = log(L_1/L_0)$$

649 Which, using Wilks's theorem, was converted to a P value for each gene. The Benjamini-

650 Hochberg procedure for controlling the False Discovery Rate (FDR) was then applied, and genes

651 with FDR <0.1 were reported as significant.

652

653 **Copy number determination by SNP array**

654 Probe-level signal intensities from Affymetrix SNP6 .CEL files for tumor samples across different

655 cancer types were combined, calibrated, normalized, and segmented in uniform fashion as

656 previously described[50]. Markers identified as having recurrent germline copy number variations

657 using normal samples were excluded. For samples for which ABSOLUTE[54] purity/ploidy calls

658 were available, copy numbers were scaled by a factor inversely proportional to the purity estimate

659 to remove the effects of admixed normal cells, as previously described[50]. Copy number profiles

660 were deconstructed into underlying somatic events using GISTIC[55]; only focal events with length

661 less than 1 megabase were used to limit normalization bias between males and females on chrX.

662 Copy number for a gene was determined by using the most extreme copy number among the

663 markers spanning the gene. The threshold for calling somatic copy number events was chosen

664 by considering the distribution of copy number across both male and female cohorts

665 (Supplementary Figure 3). Because normalization decouples chrX overall copy number from that

666 of the autosomes, retention or loss of the entire chrX was called by applying a threshold of 1.6 to

667 a robust average of the un-normalized, calibrated signal across the X chromosome

668 (Supplementary Figure 3). The maximum and minimum 0.2% probe-level signals were excluded.

669

670 **Power analysis**

671 The power calculation was performed using a binomial model similarly as previously described[39].

672 The power to detect a male-biased mutation can be determined by first calculating the maximum

673 number of mutations in males, $m_{max}$, that would be considered non-significant in the null model,

674 based on an inverse binomial cumulative distribution function, given the desired bound on the

false discovery rate (<0.1), the total number of mutations in a gene in the entire cohort, n (a function of the fraction of patients with a mutation and the total number of patients), and the fraction of all events that occur in males (as opposed to in females) in the cohort, which determines the probability that a mutation will be in a male tumor, $p_0$. The power is then determined by the probability of discovering at least that many mutations in male patients under the alternate hypothesis, which adjusts $p_0$ by the hypothesized increased relative risk of a mutation occurring in a male tumor, $R_{signal}$, the fold change compared to the null hypothesis.

$$R_{background} = p_0/(1 - p_0)$$

$$R_{total} = R_{background} * R_{signal}$$

$$p_{alt} = 1 / (1 + 1/R_{total})$$

To generate figures, we calculated the number of patients (N) required for power to detect 80% of genes with a statistically significant male bias, given hypothetical values of $R_{signal}$, n, and $p_0$.


**Y copy loss determination**

We calculated chrY coverage at a per-megabase level for paired tumor and normal sequencing data, and then normalized to total coverage of the exome. A region of chrY between bases 29000001-58000000 was omitted due to poor coverage. An additional region, 13000001-14000000, was removed due to high incidence of misalignment. A biphasic pattern of copy number quantitation was observed in male samples (Supplementary Figure 6), with a distinct population of tumors demonstrating significant median copy loss compared to the paired normal sample. ChrY loss was therefore assigned to those samples with <25% in overall coverage of chrY in the tumor compared to the paired normal sample.


**Allele-specific expression analysis**

702 Exome sequencing and RNA-seq data from somatic and germline samples from GBM, LGG,

703 HNSC, KIRC, LUSC, and LUAD sets from TCGA, and brain, lung, and whole blood sample sets

704 from GTEx were analyzed. For each, allele-specific RNA-seq pileup counts were called at

705 heterozygous germline sites and tumor sites from cancers (determined by exome sequencing)

706 from the respective analysis. Duplicate, non-primary, soft clipped and Phred quality 0 reads were

707 not included in the pileup count. For sites where at least 20 reads (tumor samples) or 8 reads

708 (normal samples) were detected by RNA-seq, the count of the less frequent allele was divided by

709 the total allele count to obtain the minor allele fraction. Data are represented as the average

710 minor allele fraction for all sites in the indicated gene, or as the average of all sites in all genes in

711 the 'non-escape' group.

712

713 **Analysis of GTEx male vs female expression data**

714 We obtained RNA-Seq expression data from the GTEx project as gene reads per sample by

715 accessing the project's pipeline (http://www.gtexportal.org/home/). For each tissue, we

716 normalized the signal across samples. The whole blood samples used were restricted to those

717 collected ante mortem. There were 85 whole blood samples (25 female and 60 male) and 360

718 brain samples (124 female and 233 male) that passed quality control and fit the selection criteria.

719 We tested for bimodality of *ATRX* expression in the brain in two ways. First, we fit a Gaussian

720 mixture model to the empirical densities using the R package "mixtools" (http://cran.r-

721 project.org/web/packages/mixtools/), function normalmixEM with parameters: k=2, epsilon = 1e-

722 08, maxit = 1000, maxrestarts=20, and assessed the distance between resulting means. Second,

723 we tested for bimodality with a likelihood ratio test for bimodality in two-component mixtures. The

724 method contrasts the likelihood of the data obtained under restricted and unrestricted maximum

725 likelihood fits of mixture of normal distributions. Under the assumption of equal variance,

726 bimodality solely depends on the mixture component weight and the ratio of the distance between

727 means and the variance. With unequal variance, it is also determined by the ratio of variances.

728  We   used   R   package   "diptest"   to   test   for   multimodality   ([http://cran.r-](http://cran.r-)

729  [project.org/web/packages/diptest/](http://cran.r-project.org/web/packages/diptest/)).

730

731  **Calculation of percentage excess male risk**

732  For a given gene in a specific disease, we calculated the excess male risk associated with loss-

733  of-function mutation of a single gene on chrX. We cannot assume that the M:F ratio in our sample

734  set is the same as the general population incidence of that disease, so we adjusted the overall

735  M:F ratio to SEER data for each cancer in the US population. The fraction of the excess male risk

736  in the disease attributable to a specific gene mutation was calculated as follows: (# males with the

737  gene mutated in our dataset – (Z * # of females with the gene mutated in our dataset)) / (# males

738  in our dataset – (Z * # of females in our dataset)); where Z = our dataset M:F ratio / SEER data

739  M:F ratio.

740

741  *Cnksr2* **knockdown, western blotting, and soft agar colony assays**

742  Mouse *Cnksr2* or control RFP shRNA constructs (oligonucleotide sequences in Supplementary

743  Table 12) were cloned per protocol ([http://www.addgene.org/tools/protocols/plko/](http://www.addgene.org/tools/protocols/plko/)) into a pLKO.1

744  vector modified to express a GFP reporter in place of the puromycin resistance gene. Lentivirus

745  was produced and murine 3T3 cells (from ATCC, mycoplasma-free) were infected using standard

746  protocols (Addgene) and GFP-positive cells were sorted after 7 days. Western blotting was

747  performed as previously described[56] using antibodies recognizing phospho-ERK (Cell Signaling

748  #4370), total ERK (Cell Signaling #9102), and tubulin (Sigma #T6074). Soft agar colony assay

749  was performed as previously described[57]. After three weeks, random microscopy fields with an

750  area of 3.29 mm$^2$ were scanned and colonies of minimum 50 um$^2$ in size were counted using a

751  CellCelector (ALS, Germany).

752

753  **RNA-sequencing and gene set enrichment analysis (GSEA) in *Cnksr2* knock-down cells**

Total RNA was prepared using a MiRNeasy kit (Qiagen). Illumnia sequencing libraries were prepared using Illumina TruSeq Stranded mRNA sample preparation kits from 500ng of purified total RNA according to the manufacturer's protocol. The finished dsDNA libraries were quantified by Qubit fluorometer, Agilent TapeStation 2200, and RT-qPCR using the Kapa Biosystems library quantification kit according to manufacturer's protocols. Uniquely indexed libraries were pooled in equimolar ratios and sequenced on an Illumina NextSeq500 with single-end 75bp reads by the Dana-Farber Cancer Institute Molecular Biology Core Facilities. Reads were aligned to the mm9 reference genome assembly using STAR (v25.1b) (https://github.com/alexdobin/STAR). FPKM expression values were calculated using cufflinks (v2.2.1) (http://cole-trapnell-lab.github.io/cufflinks/). Spearman correlation and principal component analysis were performed using VIPER (https://bitbucket.org/cfce/viper/). GSEA (http://www.broadinstitute.org/gsea/) was performed as previously described[41,58]. Network enrichment mapping was performed using Cytoscape (http://www.cytoscape.org/).

**Methods-only References**

50   Zack, T. I. *et al.* Pan-cancer patterns of somatic copy number alteration. *Nature genetics* **45**, 1134-1140, doi:10.1038/ng.2760 (2013).
51   Gao, J. *et al.* Integrative analysis of complex cancer genomics and clinical profiles using the cBioPortal. *Science signaling* **6**, pl1, doi:10.1126/scisignal.2004088 (2013).
52   Cerami, E. *et al.* The cBio cancer genomics portal: an open platform for exploring multidimensional cancer genomics data. *Cancer discovery* **2**, 401-404, doi:10.1158/2159-8290.CD-12-0095 (2012).
53   Siepel, A. *et al.* Evolutionarily conserved elements in vertebrate, insect, worm, and yeast genomes. *Genome research* **15**, 1034-1050, doi:10.1101/gr.3715005 (2005).
54   Carter, S. L. *et al.* Absolute quantification of somatic DNA alterations in human cancer. *Nature biotechnology* **30**, 413-421, doi:10.1038/nbt.2203 (2012).
55   Beroukhim, R. *et al.* Assessing the significance of chromosomal aberrations in cancer: methodology and application to glioma. *Proceedings of the National Academy of Sciences of the United States of America* **104**, 20007-20012, doi:10.1073/pnas.0710052104 (2007).
56   Yoda, A. *et al.* Mutations in G protein beta subunits promote transformation and kinase inhibitor resistance. *Nature medicine* **21**, 71-75, doi:10.1038/nm.3751 (2015).
57   Hammerman, P. S. *et al.* Mutations in the DDR2 kinase gene identify a novel therapeutic target in squamous cell lung cancer. *Cancer discovery* **1**, 78-89, doi:10.1158/2159-8274.CD-11-0005 (2011).

58    Lane, A. A. *et al.* Triplication of a 21q22 region contributes to B cell transformation through HMGN1 overexpression and loss of histone H3 Lys27 trimethylation. *Nature genetics* **46**, 618-623, doi:10.1038/ng.2949 (2014).