

Convergence of dispersed regulatory mutations predicts driver genes in prostate cancer

Richard C. Sallari^{1,2}, Nicholas A. Sinnott-Armstrong³, Juliet D. French⁴, Ken J. Kron^{5,6}, Jason Ho^{5,6}, Jason H. Moore⁷, Vuk Stambolic^{5,6}, Stacey L. Edwards⁴, Mathieu Lupien^{5,6,8} and Manolis Kellis^{1,2}

¹Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology (MIT), Cambridge, Massachusetts 02139, USA.

²The Broad Institute of MIT and Harvard, Cambridge, Massachusetts 02142, USA.

³Department of Genetics, Stanford University School of Medicine, Stanford, California 94305, USA.

⁴Cancer Division, QIMR Berghofer Medical Research Institute, Brisbane, Queensland 4029, Australia.

⁵The Princess Margaret Cancer Centre — University Health Network, Toronto, Ontario M5G 1L7, Canada.

⁶Department of Medical Biophysics, University of Toronto, Toronto, Ontario M5G 1L7, Canada.

⁷Institute for Biomedical Informatics, Perelman School of Medicine, University of Pennsylvania, Philadelphia, Pennsylvania 19104, USA.

⁸Ontario Institute for Cancer Research, Toronto, Ontario M5G 0A3, Canada.

Abstract

Cancer sequencing predicts driver genes using recurrent protein-altering mutations, but detecting recurrence for non-coding mutations remains unsolved. Here, we present a convergence framework for recurrence analysis of non-coding mutations, using three-dimensional co-localization of epigenomically-defined regions. We define the regulatory plexus of each gene as its cell-type-specific three-dimensional gene-regulatory neighborhood, inferred using Hi-C chromosomal interactions and chromatin state annotations. Using 16 matched tumor-normal prostate transcriptomes, we predict tumor-upregulated genes, and find enriched plexus mutations in distal regulatory regions normally repressed in prostate, suggesting out-of-context de-repression. Using 55 matched tumor-normal prostate genomes, we predict 15 driver genes by convergence of dispersed, low-frequency mutations into high-frequency dysregulatory events along prostate-specific plexi, controlling for mutational heterogeneity across regions, chromatin states, and patients. These play roles in growth signaling, immune evasion, mitochondrial function, and vascularization, suggesting higher-order pathway-level convergence. We experimentally validate the *PLCB4* plexus and its ability to affect the canonical PI3K cancer pathway.

Introduction

Sequencing has revealed both germ-line variants that underlie cancer risk and somatic mutations that drive cancer progression. However, disparate perspectives emerge from each. Germ-line variants identified in genome-wide association studies (GWAS) are predominantly distal to genes. When experimentally characterized they appear to interact with gene promoters by folding to their physical location in three-dimensional space. These variants have slight/subtle effects on transcription factor binding and gene expression but ultimately stack the odds towards disease development. By contrast, somatic alterations identified through tumor sequencing are far more severe and have a direct impact on gene sequences. Both emerging perspectives, the germ-line and the somatic, belong to the same disease, but it is not yet clear how they are related.

Here, we present a unifying framework to reconcile these disparate perspectives and venture into the continuum between them. In this framework we propose the existence of a theoretical genomic event in which heterogeneous variants that are scattered and far from each other on the one-dimensional genome sequence but that are physically adjacent to each other in the three-dimensional volume of the

cell nucleus manifest as a coherent cellular phenotype. We define a plexus as a set of interacting loci that are next to each other in the cell volume but scattered over the genome sequence. The number of possible plexi quickly becomes astronomical, even when they are composed of a handful of loci. Without experimental knowledge of the location of active loci and their interactions, the plexus framework is computationally and statistically intractable. Minimally, looking for a driver plexus requires whole genome sequencing of cancer-normal pairs and maps of chromatin states and chromosome interactions for the cancer's tissue of origin.

In this study, we apply the plexus framework to prostate adenocarcinoma. We use matched ChIP-seq and Hi-C to map the locations of regulatory elements and their interactions in normal prostate cells (RWPE1). We then build the prostate-specific plexus for every protein-coding gene in the human genome. We initially use the plexi to analyze dysregulated genes in 16 cancer-normal transcriptome pairs. This approach reveals that the plexi of dysregulated genes enrich in dispersed non-coding mutations that converge on gene promoters and disrupt their function. We then develop a plexus recurrence test that we apply to 55 cancer-normal whole genome pairs. This test allows us to uncover 15 driver plexi containing novel candidate cancer genes with diverse roles that further converge on growth signaling, immune evasion, mitochondrial function and vascularization. Finally, we experimentally demonstrate how our most robust result, the *PLCB4* plexus, disrupts the PI3K pathway to alter cell growth and drive tumor progression. Our results have broad implications beyond identifying cancer genes. We hope the plexus framework will boost power in sequence association studies and facilitate the interpretation of rare and private variants in the context of precision medicine.

The plexus framework

The leading paradigm for identifying driver genes in cancer genomics has been to search for recurrent mutations in multiple independent patients. This process is confounded by many factors, such as the variation in mutation rates across the genome due to transcription and replication timing¹. While cancer recurrence has typically focused on protein-coding mutations, recent studies show that regulatory regions can also be the targets of recurrent mutations. In particular, single regulatory regions linked to the *TERT*² and *TAL1* genes³. But in contrast to coding recurrence, where hundreds of genes have been identified and many more remain to be discovered⁴, the analysis of recurrence at single regulatory regions has not yielded conclusive results⁵.

Expanding the concept of recurrence to non-coding mutations poses several challenges that are currently unmet. First, cancer is a highly heterogeneous disease among patients. Because multiple regulatory loci can be associated with the same gene, each locus might be mutated in a different tumor sample, requiring methods that go beyond a single region. Second, regulatory loci can lie far from the genes they regulate. This demands precise methods for identifying interacting loci over long-range chromatin conformation loops⁶⁻⁹. Third, mutation rates are heterogeneous over the genome. Regions associated with active histone modification marks, for example, can show dramatically lower background mutation rates due to higher accessibility for the DNA repair machinery¹⁰⁻¹². Fourth, the regulatory code is poorly understood. This complicates the prioritization of mutations and the assessment of regulatory consequences^{8,13,14}. Fifth, all the parameters used in the statistical modeling vary by cell type. The interactions between loci, chromatin states, DNA accessibility and the concentrations of the transcription factors decoding the regulatory instructions are specific to each and every cell type in the human body.

Here, we directly address these challenges and introduce a theoretical and methodological framework for recurrence analysis of non-coding mutations in prostate cancer. We begin by inferring the plexus of

every protein-coding gene (**Fig. 1a**), defined as the set of all proximal and distal regulatory elements acting through intra- or inter-chromosomal interactions (also *cis* and *trans*, respectively) based on their chromatin state and their three-dimensional links to their target genes. This allows us to collapse mutations that are heterogeneous across samples and scattered over multiple genomic loci (**Fig. 1b**). Furthermore, functional annotations allow us to separate sparse driver mutations among confounding passenger mutations. Finally, we effectively aggregate mutations that are individually low in frequency into high-frequency regulatory recurrence events, based on their convergence into common target genes (**Fig. 1c**). This can be achieved even in the absence of local alterations, protein coding or otherwise.

Plexus assembly from matched Hi-C and ChIP-seq

We first establish the prostate-specific plexus of every protein-coding gene in the human genome. Regulatory annotations are highly tissue specific¹⁵. We therefore use the RWPE1 prostate cell line as a reference, and profile five histone modification marks using ChIP-Seq. We use ChromHMM¹⁶ to define eight chromatin state annotations consisting of: promoters ('pro') with strong H3K4me3 but no H3K4me1; enhancers ('enh') with H3K4me1 but no H3K4me3; regulatory elements ('reg') marked with both enhancer and promoter signatures; transcription-associated regions ('txn') with H3K36me3; poised elements ('poi') with H3K27me3 and at least one other active mark; repressed elements ('rep') with H3K27me3 only; and low-activity regions ('low') where no marks are detected (**Fig. 2a, S1**). Raw plexi exhibit an over-abundance of active regulatory states (opn, pro, reg, enh), especially through proximal interactions (**Table S1**). We treat open chromatin regions ('opn'), based on DNaseI in RWPE1¹⁷, as a separate class, regardless of their enclosing chromatin state.

We find large variation in mutation rates, across both chromatin states and tumor samples (**Fig. 2b**). Open chromatin regions show the lowest mutation rate (1.5 mutations/Mb), consistent with previous reports^{10,12,18}, attributed to their increased association with the DNA repair machinery¹¹. However, even outside DNaseI regions, mutation rates vary greatly across chromatin states (from 1.6 to 6.9 mutations/Mb on average), and across tumors (from 0.7 to 2.8 mutations/Mb for low-activity regions). Mutation rates do not correlate with GC content, CpG dinucleotide rate, nucleotide, or di-nucleotide composition (**Table S2**), and chromatin states preserve their relative mutation rates across tumors (**Fig. 2b**), suggesting sequence-independent mechanisms, possibly due to differential interactions with the repair machinery by chromatin regulators. Surprisingly, once open chromatin regions are excluded from chromatin state annotations, the expected inverse correlation between epigenomic activity and mutation rate is lost. Instead, promoter regions free of open chromatin show the second highest mutation rate.

We link these regulatory annotations to each gene using a prostate-specific map of chromosomal interactions that is also derived from the RWPE1 cell line¹⁹ (**Data S1**). With this data we generate two types of plexus for each protein-coding gene: a raw plexus, which contains any locus with evidence of interaction, and a cut plexus, in which each interaction is assessed using a permutation test and filtered based on a p-value cutoff of 0.05 (**See Methods**). While the short lengths of regulatory elements make contiguous, single-element recurrence rare (**Fig. S2**), through its plexus a gene can be associated with a much higher richness of variation across a patient cohort. The set of raw plexi associates genes with a median of 106 linked proximal elements (within 100kb), 1420 distal intra-chromosomally (*cis*), and 1824 inter-chromosomally (*trans*) interacting elements; providing an abundant source of mutation to each gene (**Fig. S3a, data S2**). Indeed, each gene interacts with a median of 21 mutations in proximal elements, 399 mutations in distal *cis*-elements, and 625 mutations in *trans*-interactions (**Fig. 2c**),

providing sufficient power to study plexus-level mutation. We use the raw plexi to identify recurrent driver events across the 55 patients; in this way we cast a broad net that we progressively tighten to further concentrate our candidates. Conversely, cut plexi associate each gene with a median of 30 interacting proximal elements (within 100kb), 319 distal *cis*-elements, and 227 *trans*-elements (**Fig. S3b, data S2**). This leads to a more focused source of mutation, with each gene being associated with a median of 1 mutation in proximal elements, 24 mutations in distal *cis*-elements, and 17 mutations in *trans*-elements (**Fig. 2d**). We exploit the stringent interactions of the cut plexi for intra-chromosomal interactions to look for enrichment of distal mutations in cancer-dysregulated genes.

Dysregulated genes enrich in plexus mutations

To study the relevance of plexus mutations as a mechanism of gene dysregulation we use whole transcriptomes obtained from 16 of the 55 patients with whole genomes²⁰. Every tumor sample has a matched normal sample from adjacent prostate tissue. From the normal tissue, we establish an expected range for every gene and use each distribution to normalize tumor transcriptomes (**Fig. 3a, see Methods**). Tumor samples showed significantly greater variance (Wilcoxon $P < 10^{-15.7}$) in expression than normal prostate samples (**Fig. 3b**). For each gene, we search for pairs of tumor samples where one gene instance is dysregulated in one patient and the other is unchanged. Pairing multiple gene instances in this manner allows us to compare a large set of dysregulated gene instances against a control set that is matched one-to-one so as to preserve the genomic and epigenomic properties between the two sets (**Fig. 3a**). Mutational properties between the two sets, which contain different compositions of patients, are corrected by incorporating the patient- and chromatin-specific mutation rates into the enrichment calculations (**See Methods**). Additionally, we only use dysregulated and unchanged gene instances when the normal prostate samples for the same two individuals show normal expression. We do this to ensure that dysregulation is not already present in the adjacent matched tissue before tumor development. Through this approach we identify 17,850 dysregulated-unchanged paired samples over a total of 2,579 genes (**Fig. 3c, data S3**), of which 83% are up regulated (14,893), and only 17% are down regulated (2,957).

We test the hypothesis that the plexi of dysregulated gene instances are enriched for mutations that are distal to the gene (>100Kb from gene body). For this we only use high confidence interactions from the cut plexi (p -value < 0.05; intra-chromosomal), and restrict our analysis to up-regulated genes where we have more power to detect an effect. We perform enrichment tests for mutations across all chromatin states and over a range of magnitudes for up-regulation. We find a consistent enrichment for 'pro' elements that increases as dysregulation becomes more extreme (Fig 3d). This enrichment peaks at 10 standard deviations from normal expression ($P_{\text{Bonf}} < 0.05$; 20,000 permutations). These enrichments remain even when we repeat the analysis removing all gene instances with copy number alterations, albeit with less statistical power (**Fig. S4**). Interestingly, 'opn' and 'reg' elements show signs of being protected from mutation; perhaps they represent a set of highly optimized activators that are less likely to increase their function through random mutation.

Based on our previous work regarding the gain and loss of enhancers in tumor initiation in colon cancer²¹ and tumor progression in breast cancer²², we hypothesize that a fraction of mutations in 'low' elements for prostate might be active in non-prostate cell lines, thus driving dysregulation activity in prostate cancer through out-of-context de-repression of existing but dormant regulatory elements, as opposed to creating them from scratch. Indeed, 'low' elements harboring mutations are strongly enriched for both promoter (Wilcoxon $P < 10^{-11}$) and enhancer states (Wilcoxon $P < 10^{-11}$) in other cell

types (**Fig. 3e**). These are active in a diverse panel of cell and tissue types²³, including immune cells, GI-tract, and ESCs (**Fig. 3f**), suggesting co-option of diverse, non-prostate elements.

Plexus recurrence test reveals hidden drivers

Having established that dysregulated genes are enriched in distal plexus mutations, we next sought to identify individual genes with an excess of mutations in their plexi in the whole genomes of the 55 tumors samples. It is highly unlikely that positive selection acts exclusively on cancer driver genes through coding and proximal promoter mutations, especially considering how most GWAS variants that increase the risk of cancer act through distal regulatory mechanism, and the bewildering diversity of mutational processes in tumor evolution.

A plexus recurrence test is faced with the same confounders as recurrence tests for coding genes, namely mutational heterogeneity across genomic regions due to cell-type-specific transcriptional activity¹ and chromatin states, in addition to variability among patient and tumor mutational signatures. In testing a plexus we must also account for changes in the confounders across the constituent loci (**Fig. 4a**). Incidentally, we find that regional mutational heterogeneity is even more extreme than previously recognized, following a power-law distribution at the 50kb scale (**Fig. 4b**), suggesting localized mutational bursts. The plexus recurrence tests starts by gathering a plexus' mutations, regional mutation rate estimates and chromatin states for all the loci it contains. These layers of information are stored at a resolution of 100bp; we refer to these intervals as 'tiles'. We retrieve a tile array for every protein-coding gene in the human genome through that gene's raw plexus. The tile array is then decomposed into the two major confounders: regional mutation rate and chromatin state (**Fig. 4c**). We compress mutation rates into 15 bins of exponentially increasing mutational intensity (taken from the tile's 50kb context). Chromatin states are stored as the 8 categories previous described.

Statistical significance is computed through permutation. The tile decomposition of a plexus is used to guide the random sampling of tiles from the whole genome so as to match the chromatin state and regional mutation rate properties of the test plexus. We then retrieve patient mutations for the permuted tile array so as to match the heterogeneity of the mutation rate in the patients. By aggregating mutations over all permutations we obtain the expected number of mutations for each patient and chromatin state over the tile array (**Fig. 4d**). The expected mutation counts allow us to convert observed mutations (**Fig. 4e**), into enrichment scores (**Fig. 4f**). Because the enrichments we previously observe for dysregulated genes in plexus mutations are highly dependent on chromatin state, we test each chromatin state separately. We combine the enrichment scores across all 55 patients to obtain a final list of p-values for each of the chromatin states and for the gene's exons (**Fig. 4g**). We refer to this procedure simply as 'the plexus recurrence test' (**See Methods**).

Applying the plexus recurrence test to the 55 prostate cancer whole-genome sequences, we identify 15 recurrently mutated plexi that are statistically significant (**Fig. 4h, table S3, data S4**). The genes varied greatly in the enriched chromatin state ('txn', 'pro', 'rep', 'poi', 'enh'), convergence rate (35%-89%), number of mutated elements (7-62), and number of mutations (24-150). The plexi do not share regulatory regions, however, the *RRAD* plexus also contains the *FAM96B* gene (**Table S4**). These 15 genes lie in a small number of common pathways, suggesting higher-order functional convergence. They are involved in cell growth, migration and proliferation through overlapping roles in androgen, insulin and circadian rhythm signaling (*INSRR*, *PLCB4*, *CRY2*, *RRAD*, *SPANX* and *SSX*), immune evasion (*ITM2A*, *IDO2*, *ZC3H12B* and *ZBED2*), mitochondrial function (*COQ3* and *SLC25A5*) and vascularization (*EDNRA*). Two genes remain uncharacterized (*C14orf180* and *ZCCHC16*). Several of these genes have been linked to cancer (*INSRR*²⁴, *RRAD*²⁵, *SSX*²⁶,) and prostate cancer, specifically

(CRY2^{27,28}). However, the most clinically relevant gene we identify is probably *IDO2*, a critical partner of *IDO1*²⁹. *IDO* genes constitute a key mechanisms of immune evasion and have recently become central targets in immuno-oncology³⁰.

***PLCB4* plexus 3C loops and epigenomic landscape**

Having identified a list of candidate driver plexi, we select one for experimental validation. Through the permutation approach we use to obtain the cut plexi, we use the p-values assigned to each of the edges of the 15 raw plexi. We then apply increasingly stringent thresholds to each plexus and recompute the plexus recurrence test at each step. As the cut plexi shrink they lose both mutated and non-mutated loci, making the recurrence signal oscillate and ultimately decay completely. The *PLCB4* plexus has the most robust recurrence signal among all 15 plexi (**Fig. 5a**). The signal comes from 79 mutations in 'txn' elements that are distal to the *PLCB4* gene body. Of these, 30 are on the same chromosome and spread over 12 loci. We perform chromatin conformation capture (3C) experiments and confirm 5 out of the 12 interactions with the *PLCB4* promoter (**Fig. 5b, S5, table S5, See Methods**). We group these elements into four loci and refer to them by their mega base coordinates on chromosome 20: 9.0 (which contains the *PLCB4* gene body), 9.2, 10.4, 18.5 and 30.0. In addition to the 3C experiments, the five loci are woven together by numerous direct and indirect Hi-C interactions (**Fig. 5c**).

The 3C-validated loci contain 15 mutations for 12 patients in the significant chromatin state, 'txn', and 36 mutations for 24 patients over all chromatin states (**Fig. 5d**). The significant ('txn') and extended (all states) mutation sets represent 22% and 44% of the 55 total patients, respectively. The recurrence frequency at the *PLCB4* plexus is high compared to coding genes identified in exome sequencing studies of prostate cancer. The top three genes by frequency, *SPOP*, *TP53* and *PTEN*, are mutated in 13%, 6% and 4% of patients, respectively³¹. However, the impact of mutations in coding regions is well understood. Annotations across 127 reference epigenomes from the Roadmap and ENCODE projects^{17,23} help infer the regulatory potential of the loci mutated in the *PLCB4* plexus. Gene dysregulation through out-of-context de-repression would require latent or poised regulatory elements to be present in the loci in cell types other than prostate. Although incomplete, the Roadmap and ENCODE collection of cell types can give some indication of regulatory activity, even when highly specific to a few non-prostate cell types. The annotations show highly diverse regulatory contexts in the ~40Kb regions containing the mutations (**Fig. 5d**). Some loci show consistent activity across all tissues, whereas others reveal striking prostate specificity. Finally, three out of the four distal loci in the *PLCB4* plexus (9.2, 10.4 and 30.0) contain eQTLs for *PLCB4* expression based on the analysis of 87 prostate samples from the Genotype-Tissue Expression (GTEx) project³² (**Fig. 5e, data S5, See Methods**).

***PLCB4* plexus disruption and the PI3K pathway**

The study of cancer recurrence in coding regions benefits from knowledge of the genetic code and allows filtering of mutations based on synonymity. In the regulatory setting we need tissue-specific annotations and models of protein-DNA binding to obtain a similar understanding. Unfortunately, a comprehensive regulatory code is still unavailable. However, we are able to infer a portion of the code in prostate cells by leveraging a subset of transcription factors binding profiles across canonical prostate cell lines. We first scanned the 15 mutations in the *PLCB4* plexus for binding of transcription factors in any human cell type using ReMap³³. Five of the seven mutations at the 30.0 locus overlap binding sites for ERG, TP63, SP1 or BRD4 (**Fig. 6a**). All of these factors are involved in gains, losses or fusions in prostate cancer³⁴⁻³⁷. We then interrogated the five canonical histone marks (H3K4me1, H3K4me3, H3K36me3, H3K27me3 and H3K27ac) and open chromatin in five additional prostate cell

lines, ranging from normal tissue (RWPE1, PWR1E) to low (LNCaP), moderate (DU145) and high (22Rv1, PC3) tumorigenicity (**Fig. 6b**). The reference for this study, RWPE1, shows active transcription marks across all four loci in addition to promoter marks at 9.2 and 30.0. The 9.2 and 10.4 loci show variability in the normal tissue and loss of histone marks across all cancer cell lines. Locus 18.5 is consistently active in all cell lines, whereas locus 30.0 shows a strong gain in enhancer marks in cell lines with high tumorigenicity and one normal cell line. These results suggest that gains and losses of activity in the *PLCB4* plexus relate to a range of cancerous phenotypes in a standard collection of prostate cell lines.

In order to assess the role of the tumor mutations in the *PLCB4* plexus on these activity patterns, we consider their effects on the binding of key transcription factors that are likely mediating the deposition of active histone marks. We built affinity models for all prostate transcription factors in the Cistrome database of ChIP-seq experiments (<http://cistrome.org/db>) using the Intra-Genomic Replicates (IGR) method⁸ (**See Methods**). We found that mutations in the 9.2 locus tend to increase the binding of AR, GR, FOXA1 and BRD4 factors, whereas mutations in the 30.0 locus tend to consistently decrease binding of BRD4, ERG, GABPA and ETV1 (**Fig. 6c**). However, the most commonly affected factor is AR, with disruptive mutations in all four loci. Each of the 14 mutations probed with IGR (single-nucleotide) had a large effect for at least one of the factors. Among the most dramatic results we find the creation of a FOXA1 binding site at locus 10.4 and the destruction of two binding sites for AR and ETV1 at the 30.0 locus (**Fig. 6d**). Interestingly, we found a significant enrichment in binding-altering mutations for ETV1 ($Q_{FDR} < 0.024$) across all 35 mutations at the *PLCB4* locus (**See Methods**).

PLCB4, or phospholipase C β 4, has been extensively studied in the context of circadian rhythms and auriculocondylar syndrome, where it has strong effects when disrupted^{38,39}. The role of *PLCB4* in prostate cancer is unknown, although it has been identified in a set of 96 genes associated with *in vivo* progression to castration-recurrent prostate cancer⁴⁰. Considering its ability to directly affect cell membrane lipid metabolism and the phosphatidylinositides, we tested the impact of *PLCB4* deletion or overexpression on the phosphoinositide-regulated PI3K/AKT signaling pathway in PC3 cells. While *PLCB4* deletion led to a decrease in PI3K/AKT signaling throughput (**Fig. 6e**), its overexpression resulted in activation of the PI3K pathway (**Fig. 6f**), revealing a direct link between the levels of *PLCB4* expression and the activity of this major oncogenic signaling cascade.

Discussion

We present the first scan for driver genes with a plexus recurrence test and, more generally, demonstrate the use of long-range chromatin loops to decipher genetic heterogeneity by collapsing the combinatorics of high-order, multi-locus interactions with a genome-wide adjacency matrix. We find that dispersed non-coding mutations that are individually too low in frequency for viable statistical analysis nevertheless converge into high frequency recurrence events. These events reveal novel driver genes with known and putative roles in prostate cancer even in the absence of proximal mutations (protein coding or otherwise), which have been the focus of previous studies. Furthermore, these genes show pathway-level convergence in androgen, insulin and circadian rhythm signaling, immune evasion, mitochondrial function and vascularization, providing new insights into biological processes known to underlie prostate cancer. Most notably, we identify *PLCB4*, which we validate as capable of affecting the canonical PI3K cancer pathway, and *IDO2*, a gene whose function is currently a central target in immuno-oncology and responsible for hundreds of millions of dollars in biotech investment and acquisitions³⁰.

We believe the plexus framework will be especially valuable for cancers with low mutation rates. Ependymoma, an extreme example, lacks any detectable mutational recurrence⁴¹. We can now qualify this statement by saying that it lacks *proximal* and contiguous mutational recurrence for any gene. However, *plexus* mutational recurrence might still be an important driver in such tumors. Indeed, many of the genes we have identified have no mutations in their gene bodies or in proximal regions. The plexus framework might also shed light into the temporal and functional interplay of diverse types of mutations through tumor initiation and progression. We hypothesize that the first somatic aberrations that propel a cell towards cancer act much like risk-associated variants. These primordial DNA mutations or epigenomic alterations are likely distal and regulatory in nature, having subtle and heterogeneous effects at first. But as these changes accumulate, triggering out-of-context de-repression of regulatory elements, they gradually converge on the hallmark pathways of cancer. Furthermore, plexus mutations are less likely to be immunogenic, allowing for the creation of a neutral evolutionary space in which pre-cancerous cells would accumulate a high degree of variability which, in turn, would provide a rich substrate for selection when it arises. Exploring this ‘plexus first’ hypothesis of cancer emergence and evolution poses a daunting challenge, as it will require the comprehensive mapping of multi-locus interactions in normal and cancerous cells across all human tissues.

Beyond cancer, the plexus framework introduced here is broadly applicable to the analysis of common, rare and private variants in any human disease or trait. Genetic heterogeneity in sequence association studies currently hinders our ability to uncover the molecular basis of heritability in complex traits. We hope that applying the plexus framework to association studies will reveal trait-associated plexi (and the genes therein) where, although each constituent locus explains only a small subset of the individuals, the whole plexus accounts for a significant proportion of the cohort. The low frequency problem of contiguous variation has led researchers of type 2 diabetes to conclude that rare variants do not contribute significantly to disease risk⁴². However, with our method the frequency of collapsed alleles increases and the number of tests decreases. Therefore, the plexus framework might boost power in such studies and allow for a more effective use of whole genome data. Similarly, the plexus framework could constitute a broader foundation for the interpretation of individual genomes. Instead of limiting analyses to the 1.5% of the genome that encodes proteins, all variants could be used. Ultimately, we wish to help advance personalized therapeutics and precision medicine by empowering those using whole-genome sequencing in the understanding, prediction and treatment of complex disease.

Author contributions. R.C.S. led the study and conceived the plexus framework in discussions with J.H.M. R.C.S., J.H.M., N.A.S.A, M.L. and M.K. developed the project. R.C.S. and N.A.S.A. designed, implemented and ran all computational analyses. S.L.E. and J.D.F. validated all interactions for the *PLCB4* plexus. V.S., M.L., J.H. and K.J.K. elucidated the role of *PLCB4* in the PI3K pathway. R.C.S., N.A.S.A., and M.K. made the figures. R.C.S. and M.K. wrote the manuscript with the help of all authors.

Acknowledgements. We thank Levi Garraway, Sylvan Baca, Cigall Kadoch, Aviad Tsherniak, Melina Claussnitzer, Angela Yen, Robert Altshuler, Pouya Kheradpour, Nezar Abdennur and Wouter Meuleman. This work was supported by NIH grants LM009012 and LM010098 to J.H.M. and R01 HG004037 to M.K., NHMRC project grant 1058415 to S.L.E. and J.D.F., Prostate Cancer Canada Movember Rising Star award RS2014-04 to ML, and National Science Foundation CAREER award 0644282 to M.K.

Correspondence. Correspondence and requests for materials should be sent to R.C.S. at sallari@mit.edu and M.K. at manoli@mit.edu.

Figure legends

Figure 1 | The plexus framework reveals the convergence of heterogeneous mutations. a, Visualization of a hypothetical plexus composed of four loci. The four loci contain a variety of active and inactive genomic elements (arcs with colored segments within an encompassing gray circle). The plexus' gene (red circle) has an active promoter ('pro', red) and a transcribed region ('txn', green). Intra- and inter-chromosomal Hi-C interactions connect the plexus' gene to regulatory elements elsewhere in the plexus (red Bezier curves). Additional Hi-C interaction help maintain the four loci together within the volume of the cell nucleus (gray Bezier curves). Dispersed mutations in active elements (putative drivers, red dots) are interspersed with mutations in inactive chromatin (likely passengers, gray dots). Active elements harboring mutations include enhancers ('enh', yellow) and poised regulatory elements ('poi', pink), among others. **b,** Dispersed active (red) and interspersed inactive (gray) mutations arranged by tumor sample. Mutational heterogeneity across samples leads to low-frequency mutation events at single elements or loci, each accounting for a common mechanism of tumorigenesis in only a small number of patients (3-5 out of 20). **c,** Aggregating low-frequency mutation events through a gene's plexus can reveal high-frequency driver events through convergent dysregulation of the same gene; a single mechanism now accounting for a majority of patient tumor samples (16 out of 20).

Figure 2 | Plexus assembly by connecting dispersed mutations to protein-coding genes. a, Chromatin states in normal prostate. We profiled five histone marks (columns) in RWPE1 cells, and used ChromHMM to learn 15 chromatin states (rows), which we further group in eight aggregate states (colors). We treat open chromatin regions, regardless of chromatin state, as a separate class (not shown). **b,** Mutation rate heterogeneity across chromatin states and tumor samples. Scatter plot showing mutation rates (x-axis) across 55 prostate tumor samples (y-axis) for eight chromatin states (colors). Tumor samples are sorted by average mutation rate in low-activity regions (low). Colored vertical bars indicate median mutation rate for each state across tumor samples. **c, d,** Linear histograms show the number of plexi (y-axis) by the mutation count in connected elements (x-axis) for raw (**c**) and cut plexi (**d**) assembled around all protein-coding genes. Plexus mutations counts are calculated separately for chromatin states (colors) and three classes of distance to the plexus' gene (columns). Colored vertical bars indicate the median number of mutations per gene for each chromatin state.

Figure 3 | Dysregulated genes are enriched in dispersed plexus mutations. a, An example of a dysregulated and unchanged gene instance pair, showing normalized expression (core score; y-axis) across 16 patients (x-axis) for exemplar gene *YJEFN3*. Tumor (top) and matched normal samples (bottom) are used to identify viable pairs across the whole transcriptome. A viable pair (red and green columns) is composed of a dysregulated (top, red, >3 SDs) and an unchanged tumor sample (top, green, <1 SD) where both samples have normal expression (bottom, dark grey, <1 SD) in the matched normal prostate tissue. Gray regions and red horizontal lines indicate the +/-1 and +/-3 SD intervals. **b,** Gene expression in tumor cells is more variable than in normal tissue. Relative standard deviation (RSD) of gene expression in 16 prostate tumor samples (y-axis) and 16 matched normal prostate samples (x-axis). Genes selected as dysregulated in viable gene instance pairs (panels **a**, **c**) are shown in red. **c,** Dysregulated and unchanged gene instance pairs are predominantly up-regulated. Scatter plot shows the normalized expression values (core score; y-axis) for every gene used in the analysis of dysregulated genes (x-axis). The most extreme dysregulated gene instance (>3std, red) and most normal unchanged instance (<1std, green) are plotted for every gene (as defined in panel **a**). All 16 normal prostate controls are shown in grey. Both up-regulated (n=2156) and down-regulated (n=423) genes are ordered by absolute expression difference between pairs (red vs. green). Vertical line

denotes the *YJEFN3* gene shown in panel **a**. **d**, Enrichment scores for mutations in the distal loci of plexi containing up-regulated genes, also referred to as mutations linked to dysregulation (high-confidence, intra-chromosomal interactions). Histograms show the \log_2 ratio of observed over expected mutation counts (y-axis) at increasing levels of dysregulation (x-axis). Significant enrichments are indicated with an asterisk ($P_{\text{Bonf}} < 0.05$; 20,000 permutations). **e**, Mutations linked to dysregulation enrich in promoters and enhancers of non-prostate tissues. Histograms showing the ratios of overlap between mutations linked to dysregulation against those linked to unchanged gene instances (fold changes; y-axis). Mutations linked to dysregulation enrich in promoters (red) and enhancers (yellow) when considering the breadth of Epigenome Roadmap cell and tissue types (x-axis). Numbers denote $-\log_{10}$ P-value of Wilcoxon enrichment. **f**, Out-of-context de-repression. Scatterplot showing the number of enhancers overlapping mutations linked to dysregulation (y-axis) against those linked to unchanged gene instances (x-axis). Both classes of mutations are contained in low-activity regions in prostate. Colors denote tissues from the Epigenome Roadmap project. Mutations linked to dysregulation are significantly more likely to lie in enhancer regions in non-prostate samples (Wilcoxon $P = 10^{-4.81}$).

Figure 4 | Plexus recurrence test identifies putative driver genes. **a**, The *ITM2A* plexus, visualized as in Fig. 1a, is used as an example. **b**, Regional mutation rate heterogeneity. Mutation rates (x-axis) assessed over a sliding window of 50Kb for every 100bp tile in the human genome (number of tiles, y-axis) follow a power-law distribution. Vertical bars denote boundaries of mutational bins used in tile resampling. Inset: Power-law distribution of regional mutation rates holds for each chromatin state. **c-f**, The plexus recurrence test for *ITM2A*. First, we tally the number of 100bp tiles in the *ITM2A* plexus and store them in a matrix structured by mutational bin (columns) and chromatin state (rows). This is the plexus tile decomposition for *ITM2A* (**c**). Second, we randomly sample 100bp tiles from the whole genome in such a way as to match the tile decomposition of the *ITM2A* plexus. From this we obtain a matrix of expected mutation counts for every patient (columns) and chromatin state (rows) combination (**d**). Third, we tally the mutations in the *ITM2A* plexus and obtain the matrix of observed mutation counts (**e**). Fourth, we use the observed and expected counts to obtain enrichment scores for each patient and chromatin state combination (**f**). Finally, the enrichment scores are aggregated across patients in order to obtain a final, permutation-based P-value for each chromatin state. Shading denotes intensity in the matrices; rows are independently normalized. **g**, The plexus recurrence test identifies 15 recurrently mutated plexi. Convergence across patients for the top plexi ranges from 35% to 89% of tumor samples (left group). When all regulatory mutations are taken into account convergence ranges from 78% to 100% of tumor samples (right group). Mutations: total number of mutations. Elements: number of mutated elements. Chromosomes: number of chromosomes containing mutated elements. Convergence element: largest percentage of patients with mutations in a single element. Convergence plexus: percentage of patients with mutations at the plexus level.

Figure 5 | The 3C structure and epigenomic landscape of the *PLCB4* plexus. **a**, Signal decay rate of the plexus recurrence test for the top 15 plexi. We re-compute p-values (y-axis) as we increase the stringency of Hi-C interactions (y-axis). The *PLCB4* plexus (gold) shows the most robust signal. **b**, Chromatin interactions between the *PLCB4* promoter and the rest of the plexus in the RWPE1 prostate epithelial cell line. Box plots show the chromatin conformation capture (3C) interaction strength between the *PLCB4* promoter and each of the other four loci in the plexus. Adjacent *EcoRI* fragments are used as controls. Error bars represent SD of three biological replicates assayed in duplicate. **c**, The 3C-validated *PLCB4* plexus visualized as in Fig. 1a with the following differences: only mutations for the significant 'txn' state are shown in red, scale marks for each locus are drawn every 5kb, and Bezier curves depict 3C-validated (red) and Hi-C (grey) interactions. **d**, Epigenome Roadmap annotations for the five loci of the 3C-validated *PLCB4* plexus (x-axis) depicted linearly across 127 cell types (y-axis);

15-state model). **e**, GTEx eQTLs for *PLCB4* in 87 prostate samples. Scatter plot depicts the $-\log_{10}$ P-value (y-axis) of *PLCB4* eQTLs for 783 variants (97, genotyped in red; 686 imputed in dark red) contained in the five of the 3C-validated *PLCB4* plexus (x-axis).

Figure 6 | *PLCB4* plexus mutations affect the binding of key prostate transcription factors.

Regulatory annotations for 15 non-coding mutations (columns) in the 3C-validated *PLCB4* plexus overlapping 'txn' state elements. **a**, Cancer transcription factor binding sites in non-prostate cells overlap mutations in the 30.0 locus. **b**, Histone marks in six prostate cell lines of increasing tumorigenic potential (from the bottom to the top) show losses at the 9.2 and 10.4 loci, but consistent activity at the 18.5 and 30.0 loci. Distance of the mutation to the closest region is indicated by saturation, where full saturation indicates direct overlap and minimum saturation indicates a 10kb distance. **c**, Intra-genomic Replicates (IGR) measures of affinity modulation for a collection of prostate related transcription factors assayed in prostate cell lines. Gains (red) and losses (blue) in binding are shown for significant results. Nested cells show results from multiple replicates of the same factor. **d**, IGR profile plots showing the most disruptive effect for each of the 14 single-nucleotide mutations in the 3C-validated *PLCB4* plexus. Each plot shows the affinity estimation (y-axis) for each of the reference (blue) and mutated (red) sequences. Affinity estimate profiles are shown over a 400bp window (x-axis) centered on the k-mer. Only the maximum affinity k-mer for each allele is presented in the final IGR result. **e**, *PLCB4* loss lowers PI3K signaling. Three independent *PLCB4*-deficient PC3 lines were engineered using CRISPR/Cas9 and their lysates immunoblotted with the indicated antibodies. **f**, *PLCB4* overexpression increases PI3K signaling. Three independent PC3 lines stably transduced with pBABE myc-*PLCB4* were isolated and their lysates immunoblotted with the indicated antibodies.

Supplementary figure legends

Supplementary Figure 1 | ChromHMM states and aggregate states. Emission parameters for the 15-state model learnt on the five core histone marks. Aggregate states represent simplified regulatory roles.

Supplementary Figure 2 | Recurrence at contiguous genomic elements is rare. **a**, Genome coverage for each of the eight aggregate chromatin states in prostate tissue (RWPE1). **b**, Number of elements (y-axis) and size of elements (x-axis) for each of the eight chromatin states, showing the distribution of regulatory element size in prostate tissue. The majority of regulatory elements are a few Kb in size. **c**, Percentage of measurements (y-axis) obtained from sliding a window of variable length (1,6 to 819 Kb) and calculating the percentage of all tumors that harbor at least one mutation in that window (x-axis). Sliding a window of 13 Kb (orange), the size of the largest regulatory elements, yields a distribution of recurrence events across all tumors that rarely exceeds 15% of prostate tumor samples.

Supplementary Figure 3 | Plexus composition across all genes. **a** and **b**. Number of genes (y-axis) and number of elements, nucleotides and mutations for raw (a) or cut (b) plexi (x-axes) for each of the eight aggregate states that are proximal (left), distal intra-chromosomal (*cis*; middle) or distal inter-chromosomal (*trans*; right) with respect to the gene body.

Supplementary figure 4 | Enrichment copy number corrected. Enrichment scores for mutations in the distal loci of plexi containing up-regulated genes, also referred to as mutations linked to dysregulation (high-confidence, intra-chromosomal interactions). Histograms show the \log_2 ratio of observed over expected mutation counts (y-axis) at increasing levels of dysregulation (x-axis).

Supplementary figure 5 | *PLCB4* 3C chromatograms. Sanger sequence chromatograms of 3C ligation products formed between the *PLCB4* promoter and mutated genomic loci.

Supplementary table legends

Supplementary Table 1 | Plexus of protein-coding genes are enriched for regulatory chromatin states. Ratio between the number of tiles observed on average for each plexus and the number of tiles expected based on genomic proportions in prostate cancer.

Supplementary Table 2 | Nucleotide and dinucleotide composition of chromatin states in prostate tissue. Nucleotide and dinucleotide counts and proportions (columns) for all eight aggregate chromatin states in prostate tissue (rows).

Supplementary Table 3 | Plexus recurrence test results for all GENCODE protein-coding genes. Expanded table of 15 cancer-associated genes with significantly mutated chromatin states in their plexus.

Supplementary Table 4 | Plexus intersections of all novel cancer-associated genes. Number of shared tiles between all pairs of genes among those identified under regulatory convergence. Only the portions of the plexus marked by the significantly mutated chromatin state are intersected.

Supplementary table 5 | *PLCB4* oligonucleotides. Nucleotide sequence of primers used in chromatin conformation capture (3C) experiments validating the *PLCB4* plexus interactions.

Supplementary datasets

Supplementary Data 1 | Number of interactions in the plexi of protein coding genes. Table header: Internal gene ID, GENCODE gene symbol, Gene class, Number of unique anchors at TSS, Number of proximal interactions, Number of distal *cis* interactions, Number of distal *trans* interactions.

Supplementary Data 2 | Number of tiles, elements and mutations in the plexi of protein-coding genes. Table header: Internal gene ID, GENCODE gene symbol, Stringency class, Annotation class, Distance class, *opn*, *pro*, *reg*, *enh*, *txn*, *poi*, *rep*, *low*.

Supplementary Data 3 | Dysregulated-unchanged gene-instance pairs. Table header: Internal gene ID, Gene normal core median, Gene normal core SD, Dysregulated index, Dysregulated patient ID, Dysregulated core score, Unchanged index, Unchanged patient ID, Unchanged core score.

Supplementary Data 4 | Plexus recurrence test results for all GENCODE protein-coding genes. *P*-values for the permutation test for regulatory convergence. Table header: Most significant *p*-value, Internal gene ID, *opn*, *pro*, *reg*, *enh*, *txn*, *poi*, *rep*, *low*, *exn*, Ensembl gene ID, Gene symbol.

Supplementary Data 5 | *PLCB4* 3C plexus eQTL results in normal prostate. Table header: Chromosome, Start, End, Gene isoform, Beta, T-statistic, P-value, FDR.

Methods

Data sources

Normal prostate Hi-C chromosome interaction data (chromatin loops) in the RWPE1 cell line were downloaded from the Gene Expression Omnibus (GEO) database at www.ncbi.nlm.nih.gov/geo (accession number: GSE37752) as part of the work of Rickman et al.¹⁹ Prostate cancer-normal whole genome and transcriptome pairs for 55 prostate adenocarcinoma patients were obtained through the database of Genotypes and Phenotypes (dbGaP) at www.ncbi.nlm.nih.gov/gap (accession number phs000447.v1.p1) and directly at the Broad Institute as part of the work of Baca et al.²⁰ Gene annotations were downloaded from GENCODE at www.gencodegenes.org (version 18). DNase annotations generated as part of the ENCODE project were downloaded from the UCSC Genome Browser at genome.ucsc.edu. Epigenome Roadmap promoter and enhancer annotations were obtained at the Broad Institute (www.broadinstitute.org/~meuleman/reg2map/HoneyBadger_release). Normal prostate transcriptomes and genotypes for eQTL analysis were obtained from the Genotype-Tissue Expression (GTEx) project at www.gtexportal.org. Additional ChIP-seq data for epigenomic and transcription factor binding analyses were downloaded from cistrome.org/db. Copy number data for the 55 prostate adenocarcinoma patients was obtained from www.cbioportal.org.

Epigenomic profiling of healthy prostate chromatin state

Healthy prostate ChIP-seq data for five core histone marks (H3K4me1, H3K4me3, H3K36me3, H3K27me3 and H3K27ac) and DNase was generated in the RWPE1 cell line with the same protocols used in Cowper-Sallari et al.⁸ except using sonication instead MNase and the Illumina HiSeq 2000 instead of the Genome Analyzer. Chromatin states were learned with ChromHMM¹⁶. We used 100bp elements in order to maintain the granularity of the smaller DNase annotations from ENCODE. The Epigenome Roadmap 15-state ChromHMM model was further aggregated into eight broader functional categories: open chromatin, promoter, regulatory (with mixed promoter and enhancer marks), enhancer, transcribed (with no other function), poised (promoter, enhancer and transcribed), repressed, and low (no marks). These aggregate states are denoted by the following three-character mnemonics: *opn*, *pro*, *reg*, *enh*, *txn*, *poi*, *rep* and *low* (Figure 1 and Supplementary Figure 1). All source data can be downloaded from www.pmgenomics.ca/lupienlab/tools.html.

Plexus assembly

Objective - The plexus framework seeks to identify cellular functions that are dysregulated through alterations distributed over multiple loci in the context of cancer recurrence and trait association studies. It is predicated on the notion that the set of loci that affect any given cellular function are likely to be non-contiguous and sparse on the one-dimensional sequence of the genome. Identifying these sets of loci from sequence alone or through exhaustive testing is currently infeasible; we therefore use experimentally derived annotations for both the locations and interactions of these sets of loci in order to address the sparseness and non-contiguity problems, respectively. Furthermore, both the locations and interactions of the loci that determine cellular functions are highly variable from one cell type to another. The use of cell types that are relevant to the trait or disease under study is critical. In principle, the framework can use any source of alteration. In this study we focus on somatic, single nucleotide variants in a cohort of 55 prostate adenocarcinoma patients, but we look forward to expanding the repertoire of alterations to germ line variants, all structural classes of mutations, and epigenetic and transcription factor binding changes. We define the plexus as the comprehensive set of genomic loci

that when altered can modulate a cellular function. In this study we focus on the expression of protein-coding genes. We first build the cell-type-specific plexus for every protein-coding gene; this includes regulatory elements that are both proximal and distal to the gene body, on both the same (*cis*) and other (*trans*) chromosomes. Specifically, we do this with histone marks and chromosome interactions (chromatin loops) obtained for the same cell-type in which the mutations originate. Ultimately, this allows us to search for genes that have more mutations in their regulatory elements or gene body than expected by chance alone. We use two types of plexus in this study: a lenient “raw” plexus, which intends to encompass as many of the true interacting loci as possible, and a stringent “cut” plexus, where interactions are filtered based on a permutations test.

Raw plexus - We retrieve transcription start site (TSS) and exon annotations from the GENCODE database (version 18) for all protein-coding genes. We segment the human genome (hg19) into 100bp tiles and assign a chromatin states to every tile (*opn, pro, reg, enh, txn, poi, rep* or *low*). Chromosome interactions (chromatin loops) were originally generated through the Hi-C sequencing technique. Interactions between loci are mediated through DNA binding proteins. Each binding site at the terminus of an interaction likely covers a few tens of bases. However, the Hi-C technique can only resolve the positions of the interaction termini to a few kilobases. This is due to the use of the HindIII restriction enzyme that cuts the DNA. The ends of Hi-C sequencing reads are unambiguously assigned to HindIII fragments, but the exact location of the terminus within the fragment cannot currently be determined. We refer to each segment of the genome contained between two HindIII restriction sites as an “anchor”. Every 100bp tile is therefore contained within a HindIII anchor, and anchors are connected to each other through Hi-C interactions. Through this network we assign tiles to the plexus of every gene. First we retrieve all anchors within 10Kb up and downstream of the gene’s TSS. For each anchor at the TSS, we then retrieve all other anchors connected through Hi-C interactions. Each of these anchors is extended 10Kb in both directions. We then store all tiles that overlap either a TSS anchor, any of the gene’s exons, or any of the extended anchors at the distal ends of the Hi-C interactions. The list of 100bp tiles is then sorted and filtered down to an array of unique tile indices. Raw plexus tile arrays frequently span multiple chromosomes and several megabases (**Fig. S3**). They also contain a variety of chromatin states and disparate regions of the genome with highly discordant mutation rates. The plexus tile array is the representation of all the proximal and distal elements that potentially impinge on a gene’s function. We compute plexus tile arrays for every protein-coding gene in the human genome (**Data S1**), and then retrieve tumor mutations for every array (**Data S2**). The raw plexi constitute the basis for our plexus recurrence test.

Cut plexus - Several factors can account for the presence of reads spanning two loci in a Hi-C library, many of which are confounders to identifying true regulatory interactions. We set out to assign a measure of confidence on putative interactions present in each of the raw plexi previously defined. The cut plexus program takes the raw plexus of a gene as input. It then loads the Hi-C matrices for the chromosome that contains the gene. The matrices are binned at 10Kb intervals for the intra-chromosomal interactions (inter-chromosomal are discarded for the cut plexi) and computed using the hiclib library (mirnylab.bitbucket.org/hiclib). The raw data from Rickman et al. contains four replicates 'GFP1', 'GFP2', 'ERG1' and 'ERG2' of varying sequencing depth. Interactions are tested separately across the four replicates and aggregated at the very end of execution. The GFP and ERG conditions are intended to simulate normal and cancerous cell states. We want to capture interactions that can occur across the carcinogenic continuum. Furthermore, we are primarily interested in identifying loci that can interact rather than loci that interact frequently. Therefore, because these loops might vary over time, identifying a loop in a single replicate is valid. We then assess each of the locus pairs in the

raw plexus. Interactions that are less than 100Kb from the gene promoter are tagged as “proximal”. Because for now we are only concerned with intra-chromosomal interactions, the rest are tagged as “distal_cis”. We generate a p-value for each of the four replicates by permuting matrix counts between the two loci.

Hi-C null - This null distribution represents the expected counts given the distance between loci and the genomic background at each of the loci. To compute the null we start by tallying the observed matrix count. This number is the largest value among the cells corresponding to the bins containing the pair of loci or any of the eight adjacent cells. We do this in order to include arrangements in which the probed loci and interaction anchors are separated in adjacent bins. The anchors that mediate the interactions between genes and regulatory elements might not overlap perfectly⁴³. True interactions can be mediated by physical interactions that are some kilobases away. Additionally, restriction fragment processing makes the exact location of these interactions uncertain. We then take a segment of the diagonal that crosses the cell addressed by bin coordinates of the loci. This diagonal slice of the matrix expands 100 cells upstream and downstream parallel to the diagonal and two cells perpendicular. We tally the matrix values (balanced reads) over the diagonal slice. And then randomly reassign the same number of reads over a matrix of the same dimensions as the diagonal slice. We count the number of times that we see values larger than or equal to the observed matrix count in the permuted matrix. We then regenerate the randomized matrix 100 times to improve this estimate. The p-value for each replicate is calculated by dividing the number of cells with tallies larger or equal than the observed value over the total number of cells and matrix randomizations. This is equivalent to averaging the expectation estimate over multiple permutations. We combine the p-values over the four replicates using Fisher’s method. This final p-value is what we use to filter the locus pairs in the raw plexus to generate cut plexi of varying stringency. Through this approach we attempt to estimate the probability of the observed read count given the genomic context of the two loci. Our intention is to correct for the interface or area of interaction between the loci at the level of chromosome territories. We do this at a scale that is much larger than the individual loci but that retains the local conformational background. We run the permutation test on all GENCODE V18 genes marked as ‘protein_coding’ for all chromosomes except Y and M, which yields a total of 20,233 genes. We use a p-value threshold of 0.05 to define all subsequent statistics for cut plexi.

Dysregulated gene analysis

Core scores - We identify dysregulated copies of genes in specific patients, referred to as “gene instances”, across 16 cancer transcriptomes by normalizing every gene’s expression using an aggregate of values in the matched normal prostate tissue. The normalized measure we use is a slight variation on z-scores, which we refer to as “core scores”. We calculate the median and standard deviation (SD) of each gene from the 16 normal transcriptomes. Because prostate cancer has a strong genetic component we consider that some gene instances might already be dysregulated in the matched normal tissue. Therefore, we discard the six most extreme values for each gene. We do this by sorting all gene-instances for each gene and finding the sequence of ten consecutive instances with the smallest SD. Because removing the extremes of any distribution will warp the estimates of its SD, we estimate the distortion factor extracting the core from random sets of 16 values sampled from the normal distribution and correct the core-scores accordingly (distortion factor ~2). Finally, we filter out lowly expressed genes where the normal core has a median FPKM value under 0.3{Ramskold:2009wu}.

Matching of pairs - We search for genes where at least one instance has an absolute core score larger than three SDs (dysregulated) and one other instance has an absolute core score smaller than one SD

(unchanged). These dysregulated and unchanged gene instance pairs allow us to study transcriptional dysregulation in one patient while having a control instance of the same gene in a second patient. The expression values for patients in normal tissue need to have an absolute core score smaller than one SD in order to be included in the analysis. We do this to ensure that dysregulation is not already present in the adjacent matched tissue before tumor development. We identify 17,850 viable, dysregulated and unchanged gene instance pairs over a total of 2,579 genes (**Fig. 3c, data S3**), of which 83% are up regulated (14,893), and only 17% are down regulated (2,957).

Enrichment scores - We use the cut plexi (p -value < 0.05) to assign mutations to each of the gene instances for all viable pairs. We use cut plexi and intra-chromosomal interactions in order to enrich for true interactions, as the Hi-C is noisy. Because we are strictly interested in the effect of distal mutations, we only consider those that are beyond 100Kb of the plexus' gene promoter. Additionally, we restrict our analysis to up-regulated genes where we have more power to detect an effect. Having gathered mutations for both dysregulated and unchanged gene instances in the tumor samples, we can compute the enrichment of mutations in dysregulated gene instances by comparing the two groups in aggregate. Because the specific contribution of different classes of regulatory elements to tumorigenesis is unknown, we compute enrichment scores for each of the chromatin states separately.

Pair resampling - We test several intensities of dysregulation for enrichment in mutations, from core scores between 3 and 12 SDs. For each baseline we pick all the pairs in the 17,850 viable pair set where the core score of the dysregulated gene instance exceeds or equals the baseline core score. A permuted set of pairs of the same size as the one defined by the core score baseline is resampled within the same collection of genes by randomizing the indices of the dysregulated and unchanged instances between genes. This is similar to randomly sampling 16x16 patient sample pairs, but with a non-uniform distribution over the patients. By simply permuting the indices between genes, the distribution approximately preserves the patient proportions among the dysregulated and unchanged categories across all genes. We then retrieve the plexus mutations for the dysregulated and unchanged gene instances and append them to their respective mutation matrices (gene instances by chromatin states). As we increase the core score baseline, the number of viable pairs decreases and the analysis loses power.

Patient balancing - Patient proportions in the dysregulated and unchanged matrices are not balanced. When selecting gene instance pairs for a given core score baseline, we tally the number of times each patient appears in either category. Based on the mutation rate heterogeneity across patients, we compute the expected dysregulated to unchanged ratio of mutations for each chromatin state. We then compare this to the dysregulated to unchanged ratios between the two mutation matrices. The ratio of ratios is \log_2 transformed and constitutes the enrichment score for each chromatin state at each baseline. We derive a p -value for this test by counting the number of times that the observed enrichment score is larger than the enrichment scores in the permuted gene instance pairs. Finally, we correct the p -values by the 10 core scores baselines and 8 chromatin states tested (Bonferroni; 80 tests).

Copy number correction - If dysregulation is due to regional copy number alteration, then the region is likely to have more mutations assigned to it during variant calling, as the copies are conflated in the variant calling. This could lead to the appearance of enrichment when compared to an unchanged gene instances in a patient with no copy number alterations in that region. We therefore add an additional condition to the previous approach in order to avoid the possible effect of copy number alterations in our study of dysregulated genes. When collecting gene instance pairs for each core score baseline, a pair is only considered if it has copy number data available and neither gene instance has a copy number alteration. This filter is applied to both the observed and permuted pair sets.

Plexus recurrence test

Tile resampling - The plexus recurrence test is designed to identify gene plexi that harbor more mutations than are expected by chance alone. It estimates the expected number of mutations through resampling. Its null distribution accounts for critical confounders that have been previously identified in the search for driver events in cancer genomes¹. Namely, it accounts for heterogeneous mutation rates across patients, chromatin states and genomic regions. The null distribution is computed from a resampling matrix in which all 100bp tiles used to assemble the plexi are binned based on their chromatin state and regional mutation rate. The regional mutation rate at each tile is calculated from its surrounding 50Kb by pooling mutations across all patients. We find that larger windows fail to account for the broad variations in mutation rates, whereas smaller windows are similar in size to genes, which can be legitimate units of selection. Each tile is assigned its mutational bin number based on the \log_2 of its regional mutation rate (**Fig. 4b**). The two-dimensional binning of all tiles in the human genome results in a matrix of tile arrays of 15 mutation bins by 8 chromatin states. Each tile retains the assignment of mutations for each patient; thus preserving mutation rate heterogeneity within the cohort. A similar two-dimensional binning is applied to the tile array of the test plexus, where instead of tile index list we store tile tallies (**Fig. 4c**). We refer to this matrix of tile tallies as the “tile decomposition” of the plexus, it encodes our understanding of the attributes across the loci that compose the plexus that affect its accumulation of mutations but that are independent of selection; i.e. the mutational confounders. The plexus tile decomposition guides the resampling in each permutation. Tiles are picked at random and with replacement from the resampling matrix such as to match the tile decomposition of the test plexus and control for the mutational confounders. The collection of permuted tile arrays constitutes the null distribution with which to estimate expected mutation rates for each patient and chromatin state combination.

Centroid score - In the next step, the plexus recurrence test retrieves the patient mutations for the test plexus' and each of the permutation's tile arrays. A tile is considered mutated for a patient if it contains one or more mutations. Each tile is therefore associated with a binary array with an entry for each patient. Mutations for the permuted tile arrays retain their original assignments to patients in order to control for mutation rate heterogeneity in the cohort, including patient-specific variation in chromatin state mutation rates. Observed mutations in the test plexus' tile array are compared to mutations in the permutation's tile arrays using a separate centroid score for each chromatin state and the exon annotations contained in the plexus (referred to just as “states” for brevity). The centroid score is computed in the following manner. Mutations are tallied in state by patient matrices for all tile arrays (**Fig. 4d and 4e**). Tally matrices for the permutations are stored as a three-dimensional null volume. Each list of mutation tallies for each state and patient combination is upper quartile normalized across the permutation dimension. Positive normalized tallies summed across patients constitute the centroid score, which are calculated separately for each state. This produces nine centroid scores for the test plexus and each of the permutations. Each patient contributes to the centroid score proportionally to how much it positively deviates from what is expected for that patient and chromatin state combination (**Figure 4f**).

Significance - Statistical significance for the test plexus is computed by comparing its centroid scores against those of the permutations. The size of the null distribution (number of permutations) is increased until a reliable *P*-value is determined. Every plexus is tested with at least one thousand permutations. The number increases by an order of ten in each subsequent round. Only plexi in which one or more of the nine tests (eight chromatin states and exons) have two or fewer permutations that

have a centroid score larger or equal to that of the null pass on to the next round. Obtaining *P*-values for all plexi requires up to ten million permutations in some cases, which is consistent with the number of hypotheses tested. We perform nine tests on each of the 20,318 protein-coding genes and correct the *P*-values accordingly with the Benjamini-Hochberg method (**Supplementary data 4**).

Discussion - Many features remain to be incorporated and explored in future version of the plexus recurrence test. First, due to limited availability and high cost of Hi-C data, we constructed the plexus of each gene in reference to a single prostate cell line. This single reference practice is common in functional genomic studies of GWAS trait-associated variants, which rely on small numbers of model cell lines relevant to the trait or disease being studied. However, profiling chromatin states, DNase hypersensitivity regions, and especially Hi-C interactions in each tumor and normal sample individually would provide for a much richer analysis of each patient's disease. The direct incorporation of variability in plexus structure between individuals and the regulatory rewiring within tumors would be particularly interesting. Second, even though we corrected for mutation rate differences between chromatin states, we treated all mutations in the same chromatin state as equally likely to have a regulatory effect and drive tumorigenesis. Protein-coding models that distinguish between synonymous and non-synonymous mutations and attempt to predict the effect of variants are readily available. However, similar models for regulatory alterations are not as developed. Richer regulatory models will allow future iterations of the test to incorporate the magnitude and direction of effect for non-coding mutations in the expectations derived from resampling. For example, by pre-computing the likelihood of a mutation perturbing enhancer activity or modulating binding of a transcription factor. The models will eventually leverage tumor-specific information on regulator activity, such as the intra-cellular concentration of aberrant transcription factors. Third, the three-dimensional structure of the genome has been shown to be scale-free, with rich patterns of chromosome interactions at multiple scales. Therefore, our test of non-contiguous recurrence should not be limited to a single scale, in this case the organization of regulatory elements around a single gene, but should consider both smaller and larger structures. The tested units could range from enhancer clusters or sets of interaction termini to transcriptional factories or topological domains. The pathway-level convergence we observe among the 15 drivers we identify suggests another set of layered recurrence tests across the hierarchy of cellular functions. A hierarchical recurrence test could reveal convergent mutations in the merged plexi of protein complexes, metabolic pathways and cancer hallmarks. Finally, the plexus recurrence test should combine somatic and germ line variants of strong and weak effects, that are rare to common in frequency, both protein-coding and non-coding, to cover the entire causal timeline of cancer, from inherent risk to emergence and evolution. This would constitute the first approximation of a unified model of tumorigenesis

PLCB4 plexus chromatin conformation capture (3C)

The RWPE-1 (normal prostate epithelial) cell line was kindly provided by Dr Jyotsna Batra (Queensland University of Technology, Brisbane, Australia) and cultured in KSFM supplemented with 5 ng/ml epidermal growth factor, 25 μ g/ml bovine pituitary extract and 2 mM glutamine. 3C libraries were generated using *EcoRI* as described previously{Ghoussaini:2014uq}. 3C interactions were quantitated by real-time PCR (Q-PCR) using primers designed within *EcoRI* restriction fragments (**Fig. S5, table S5**). Q-PCR was performed on a RotorGene 6000 using MyTaq HS DNA polymerase (Bioline) with the addition of 5 mM of Syto9, annealing temperature of 66°C and extension of 30sec. 3C analyses were performed in three independent library preparations with each experiment quantified in duplicate. Bacterial artificial chromosome (BAC) clones were used to create artificial libraries of ligation products

in order to normalize for PCR efficiency. Q-PCR products were electrophoresed on 2% agarose gels, gel purified and sequenced to verify the 3C product.

eQTL analysis of variants in the *PLCB4* plexus

The Genotype-Tissue Expression (GTEx) project is a publicly available resource that provides matched genotype and expression data from normal human donors{Ardlie:2015tv}. 87 of these donors have transcriptome data (RNA-seq) for prostate, in addition to genotype and covariates data. We intersect the five 3C-validated loci with the GTEx genotypes and obtain 783 variants (97 genotyped, 686 imputed) for the *PLCB4* plexus. Nine transcripts of *PLCB4* (ENSG00000101333.12) have detectable expression in at least one of the 87 prostate transcriptomes (ENST00000278655.4; a, ENST00000334005.3; b, ENST00000378473.3; c, ENST00000378501.2; d, ENST00000416836.1; e, ENST00000464199.1; f, ENST00000473151.1; g, ENST00000482123.1; h, ENST00000492632.1; i). We compute expression quantitative trait loci (eQTLs) for the nine transcripts and the 783 variants using the Matrix eQTL R package⁴⁴. The Bonferroni cutoff for statistical significance is $10^{-5.149}$ ($0.05 / (783 * 9)$). We used the age, race and ethnicity of the 87 GTEx donors as covariates in the analysis.

Intra-genomic replicates (IGR) analysis of *PLCB4* plexus mutations

We downloaded signal tracks from Cistrome DB for 118 ChIP-Seq experiments performed in prostate cell lines, as well as ERG and GABPA in Jurkat cells⁴⁵. We construct 400bp window 7mer and 8mer IGR models for each track, as previously described⁸ except that we do not apply the competition filter. We use 373,359 active prostate elements genome wide as the IGR regional filter. We define these regions as the union of peaks across 29 experiments profiling DNase, H3K27ac, H3K4me1, H3K4me2, and H3K4me3 in the same prostate cell lines. Using these models, we run the IGR algorithm on 14 out of the 15 mutations in 'txn' states and 35 of the 36 total mutations contained in the *PLCB4* plexus (only single-nucleotide variants). Many of these mutations show statistically significant affinity modulation of transcription factor binding between the reference and alternate alleles (Bonferroni corrected).

IGR filters - We have updated the original IGR program with two new filters. First, in order to discard noisy affinity models lacking sufficient instances of a given k-mer genome wide to make a clean prediction we devise the "quality" and "symmetry" filters. IGR computes an averaged binding profile for every k-mer; in this case along a 400bp window centred on the k-mer. Every k-mer has two profiles for its forward and reverse complement orientations, and every IGR result has two final k-mers for the highest affinity among the reference and alternate allele k-mer sets. The correlation between the forward profile and the mirror image of the reverse profile constitutes the measure of quality. The correlation between the forward profile and its own mirror image constitutes the measure of symmetry. We then remove any mutations results where either reference and alternate final k-mer profiles had either symmetry or quality smaller than 0.5 and both had symmetry and quality smaller than 0.85. Second, in order to only select mutations for which the effect size was large enough we calculated baseline-offset affinities for each of the final k-mer profiles. This measure compares the affinity centred at the k-mer, minus the average of the signal 195-200bp away from the k-mer in both forward and reverse orientations. Using these, we define the "maximum prominence" as the highest absolute baseline-corrected affinity in either the reference or alternate allele within 200bp of the k-mer and the "maximum difference" as the largest absolute difference between baseline-corrected reference and alternate alleles within 200bp of the k-mer. We exclude all mutation results for which the ratio between the maximum difference and maximum prominence was less than 0.5.

Enrichment analysis - We test *PLCB4* plexus mutations for enrichment of affinity modulating results that satisfy all of the previous filters. We assess the set of mutations in the *PLCB4* plexus as a whole using

Fisher's exact test within each experimental setup. Enrichments were corrected using FDR and only significant IGR mutations were used when counting (q -value < 0.024; estimate = 6.73 for 8-mers)

PLCB4 overexpression

PC3 cells were transfected with PolyJet (SignaGen Laboratories) as per manufacturer's instructions. Briefly, cells were plated the night before transfection to achieve 70% confluence on 10 cm dishes. PC3 cells were transfected with empty vector, pcDNA3.1-mycHis-PLCB4, p3xFlag-CMV10-PTEN or both. Forty-eight hours later, the cells were harvested by washing and scraping in ice-cold PBS followed by centrifugation at 1500 x g for 5 minutes at 4°C. The cell pellet was lysed in 160 μ L of CHAPS lysis buffer (40 mM HEPES pH 7.5, 0.3% CHAPS, 120 mM NaCl, 1 mM EDTA, 50 mM sodium fluoride, 20 mM beta-glycerophosphate, protease inhibitor cocktail) on ice for 20 minutes. The cell lysate was clarified by centrifugation at 15,000g for 15 minutes at 4°C and the supernatant was normalized for total protein using the Bradford assay (Biorad). SDS loading buffer was added to the normalized samples and 30 μ g of total protein was loaded on 8% acrylamide gels and transferred to PVDF membranes. The membranes were blocked in 5% BSA in TBST and immunoblotted with anti-PLCB4 (sc-20760, Santa Cruz Biotech), anti-phospho Akt (#4058, Cell Signaling), anti-phospho Erk (#9106, Cell Signaling), anti-Akt (#4691, Cell Signaling), anti-Erk (#9102, Cell Signaling) and anti-PTEN (#9559, Cell Signaling).

References

1. Mutational heterogeneity in cancer and the search for new cancer-associated genes. **499**, 214–218 (2013).
2. Highly recurrent TERT promoter mutations in human melanoma. **339**, 957–959 (2013).
3. Patel, B. *et al.* Aberrant TAL1 activation is mediated by an interchromosomal interaction in human T-cell acute lymphoblastic leukemia. *Leukemia* **28**, 349–361 (2014).
4. Lawrence, M. S. *et al.* Discovery and saturation analysis of cancer genes across 21 tumour types. *Nature* **505**, 495–501 (2014).
5. Nik-Zainal, S. *et al.* Landscape of somatic mutations in 560 breast cancer whole-genome sequences. *Nature* **534**, 47–54 (2016).
6. 8q24 prostate, breast, and colon cancer risk loci show tissue-specific long-range interaction with MYC. **107**, 9742–9746 (2010).
7. Integrative functional genomics identifies an enhancer looping to the SOX9 gene disrupted by the 17q24.3 prostate cancer risk locus. **22**, 1437–1446 (2012).
8. Breast cancer risk-associated SNPs modulate the affinity of chromatin for FOXA1 and alter gene expression. **44**, 1191–1198 (2012).
9. Pomerantz, M. M. *et al.* The 8q24 cancer risk variant rs6983267 shows long-range interaction with MYC in colorectal cancer. *Nat. Genet.* **41**, 882–884 (2009).
10. Human mutation rate associated with DNA replication timing. **41**, 393–395 (2009).
11. Liu, L., De, S. & Michor, F. DNA replication timing and higher-order nuclear organization determine single-nucleotide substitution patterns in cancer genomes. *Nat Comms* **4**, 1502 (2013).
12. Reduced local mutation density in regulatory DNA of cancer genomes is linked to DNA repair. **32**, 71–75 (2013).
13. Ernst, J. & Kellis, M. Interplay between chromatin state, regulator binding, and regulatory motifs in six human cell types. **23**, 1142–1154 (2013).
14. Khurana, E. *et al.* Integrative annotation of variants from 1092 humans: application to cancer genomics. *Science* **342**, 1235587 (2013).

15. Lupien, M. *et al.* FoxA1 translates epigenetic signatures into enhancer-driven lineage-specific transcription. **132**, 958–970 (2008).
16. Ernst, J. & Kellis, M. ChromHMM: automating chromatin-state discovery and characterization. *Nat. Methods* **9**, 215–216 (2012).
17. Consortium, T. E. P. An integrated encyclopedia of DNA elements in the human genome. **489**, 57–74 (2012).
18. Polak, P., Querfurth, R. & Arndt, P. F. The evolution of transcription-associated biases of mutations across vertebrates. *BMC Evolutionary Biology* **10**, 187 (2010).
19. Rickman, D. S. *et al.* Oncogene-mediated alterations in chromatin conformation. **109**, 9083–9088 (2012).
20. Punctuated Evolution of Prostate Cancer Genomes. **153**, 666–677 (2013).
21. Akhtar-Zaidi, B. *et al.* Epigenomic enhancer profiling defines a signature of colon cancer. **336**, 736–739 (2012).
22. Magnani, L. *et al.* Genome-wide reprogramming of the chromatin landscape underlies endocrine therapy resistance in breast cancer. **110**, E1490–9 (2013).
23. Roadmap Epigenomics Consortium *et al.* EC00: Integrative analysis of 111 reference human epigenomes. *In revisions* (2014).
24. Hua, X. *et al.* DrGaP: A Powerful Tool for Identifying Driver Genes and Pathways in Cancer Sequencing Studies. *The American Journal of Human Genetics* **93**, 439–451 (2013).
25. Reynet, C. & Kahn, C. R. Rad: a member of the Ras family overexpressed in muscle of type II diabetic humans. *Science* **262**, 1441–1444 (1993).
26. D'Arcy, P., Maruwge, W., Wolahan, B., Ma, L. & Brodin, B. Oncogenic Functions of the Cancer-Testis Antigen SSX on the Proliferation, Survival, and Signaling Pathways of Cancer Cells. **9**, e95136 (2014).
27. Chu, L. W. *et al.* Variants in circadian genes and prostate cancer risk: a population-based study in China. *Prostate Cancer Prostatic Dis* **11**, 342–348 (2007).
28. Zhu, Y. *et al.* Testing the circadian gene hypothesis in prostate cancer: a population-based case-control study. *Cancer Research* **69**, 9315–9322 (2009).
29. Metz, R. *et al.* IDO2 is critical for IDO1-mediated T-cell regulation and exerts a non-redundant function in inflammation. *Int. Immunol.* **26**, 357–367 (2014).
30. Sheridan, C. IDO inhibitors move center stage in immuno-oncology. *Nature biotechnology* **33**, 321–322 (2015).
31. Barbieri, C. E. *et al.* Exome sequencing identifies recurrent SPOP, FOXA1 and MED12 mutations in prostate cancer. *Nat. Genet.* **44**, 685–689 (2012).
32. GTEx Consortium. Human genomics. The Genotype-Tissue Expression (GTEx) pilot analysis: multitissue gene regulation in humans. *Science* **348**, 648–660 (2015).
33. Griffon, A. *et al.* Integrative analysis of public ChIP-seq experiments reveals a complex multi-cell regulatory landscape. *Nucleic Acids Res.* **43**, e27–e27 (2015).
34. Sankpal, U. T., Goodison, S., Abdelrahim, M. & Basha, R. Targeting Sp1 transcription factors in prostate cancer therapy. *Med Chem* **7**, 518–525 (2011).
35. Tucci, P. *et al.* Loss of p63 and its microRNA-205 target results in enhanced cell migration and metastasis in prostate cancer. *Proc. Natl. Acad. Sci. U.S.A.* **109**, 15312–15317 (2012).
36. Tomlins, S. A. *et al.* Recurrent fusion of TMPRSS2 and ETS transcription factor genes in prostate cancer. *Science* **310**, 644–648 (2005).
37. Asangani, I. A. *et al.* Therapeutic targeting of BET bromodomain proteins in castration-resistant prostate cancer. *Nature* **510**, 278–282 (2014).
38. Park, D. *et al.* Translation of clock rhythmicity into neural firing in suprachiasmatic nucleus

- requires mGluR-PLCbeta4 signaling. *Nat. Neurosci.* **6**, 337–338 (2003).
39. Rieder, M. J. *et al.* A human homeotic transformation resulting from mutations in *PLCB4* and *GNAI3* causes auriculocondylar syndrome. *American Journal of Human Genetics* **90**, 907–914 (2012).
 40. Romanuik, T. L. *et al.* LNCaP Atlas: gene expression associated with in vivo progression to castration-recurrent prostate cancer. *BMC Med Genomics* **3**, 43 (2010).
 41. Epigenomic alterations define lethal CIMP-positive ependymomas of infancy. **506**, 445–450 (2014).
 42. Fuchsberger, C. *et al.* The genetic architecture of type 2 diabetes. *Nature* **536**, 41–47 (2016).
 43. Bailey, S. D. *et al.* ZNF143 provides sequence specificity to secure chromatin interactions at gene promoters. *Nat Comms* **2**, 6186 (2015).
 44. Shabalin, A. A. Matrix eQTL: ultra fast eQTL analysis via large matrix operations. *Bioinformatics* **28**, 1353–1358 (2012).
 45. Sharma, N. L. *et al.* The ETS family member GABP α modulates androgen receptor signalling and mediates an aggressive phenotype in prostate cancer. *Nucleic Acids Res.* **42**, 6256–6269 (2014).

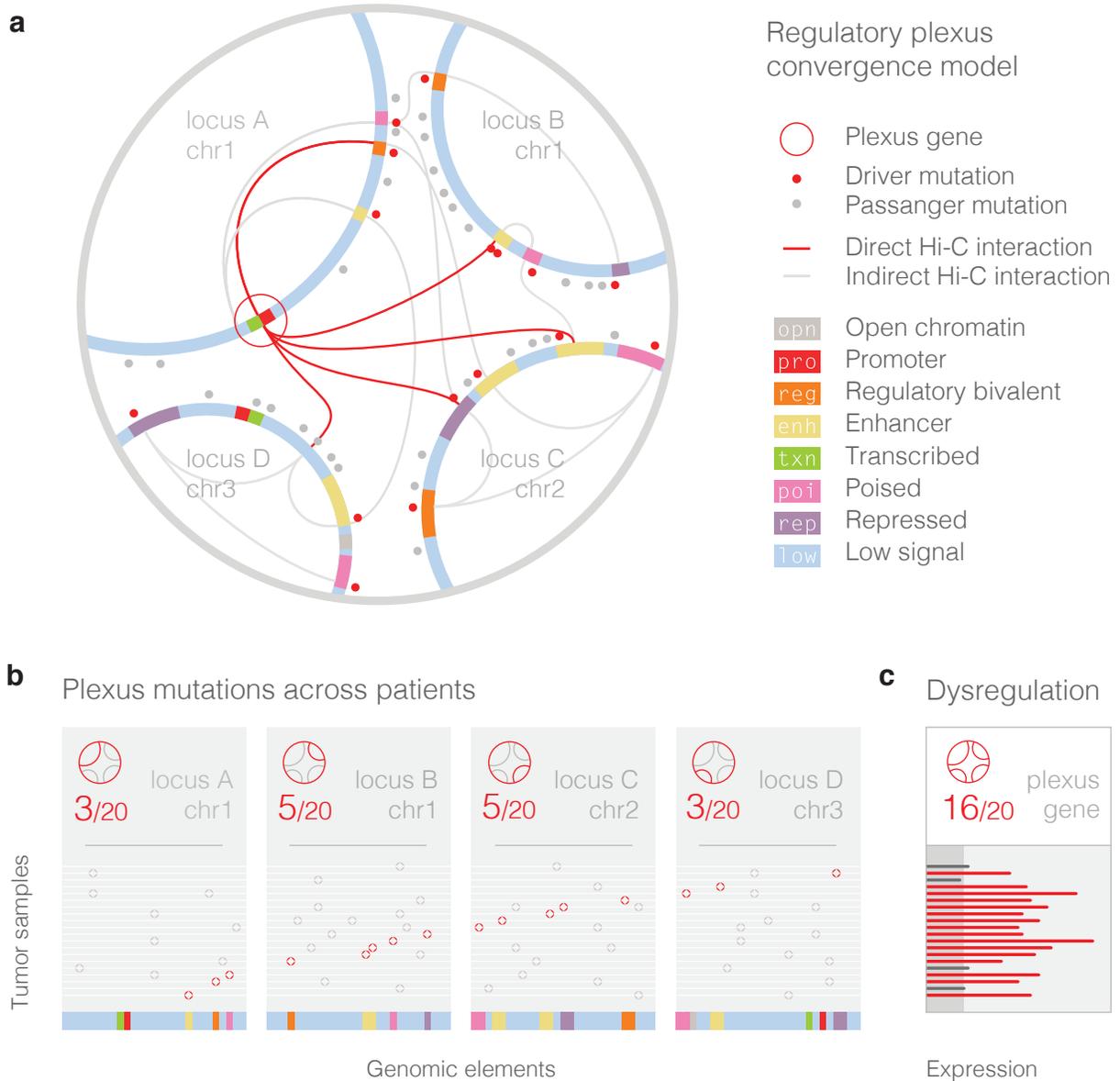


Figure 1

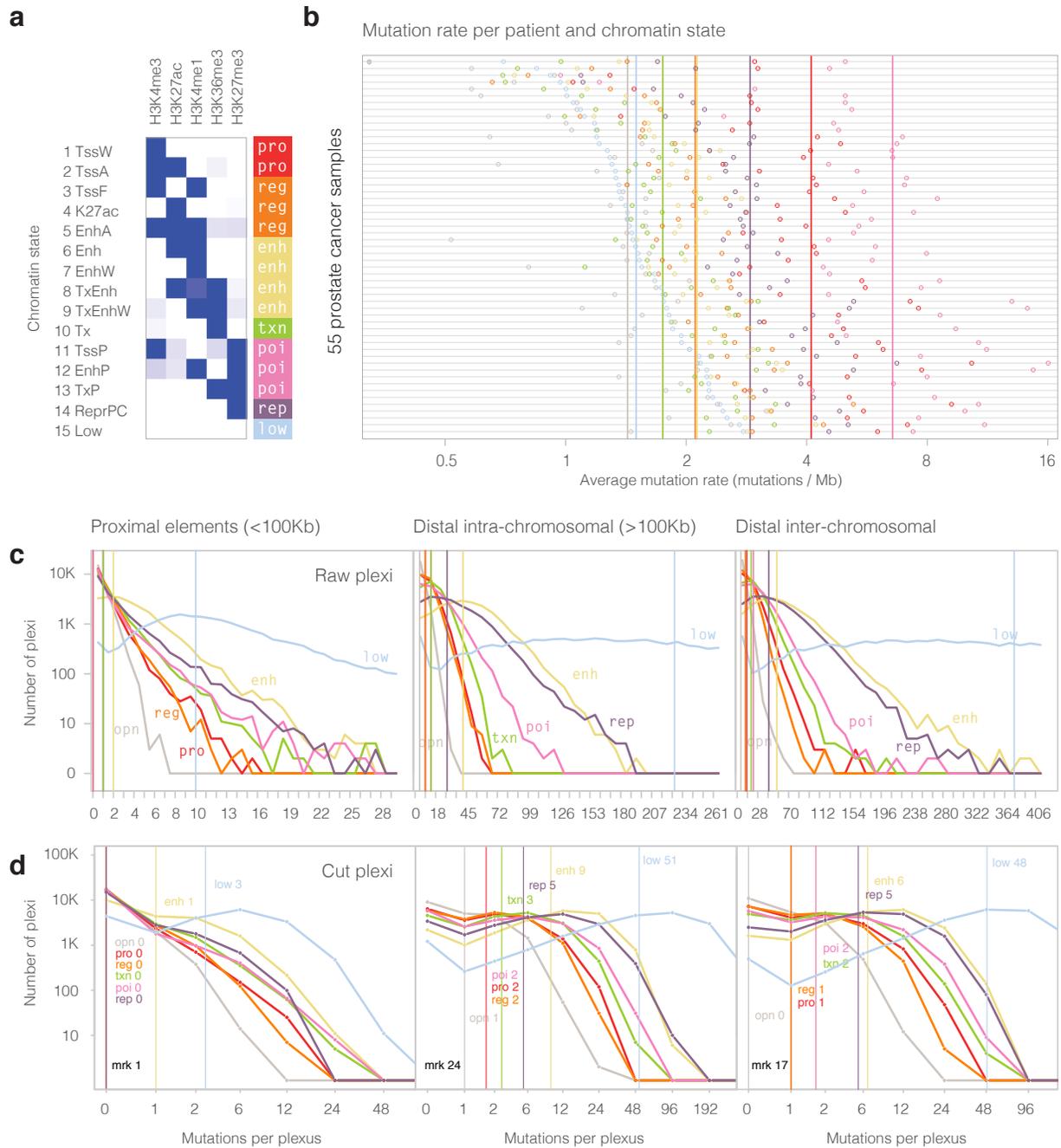


Figure 2

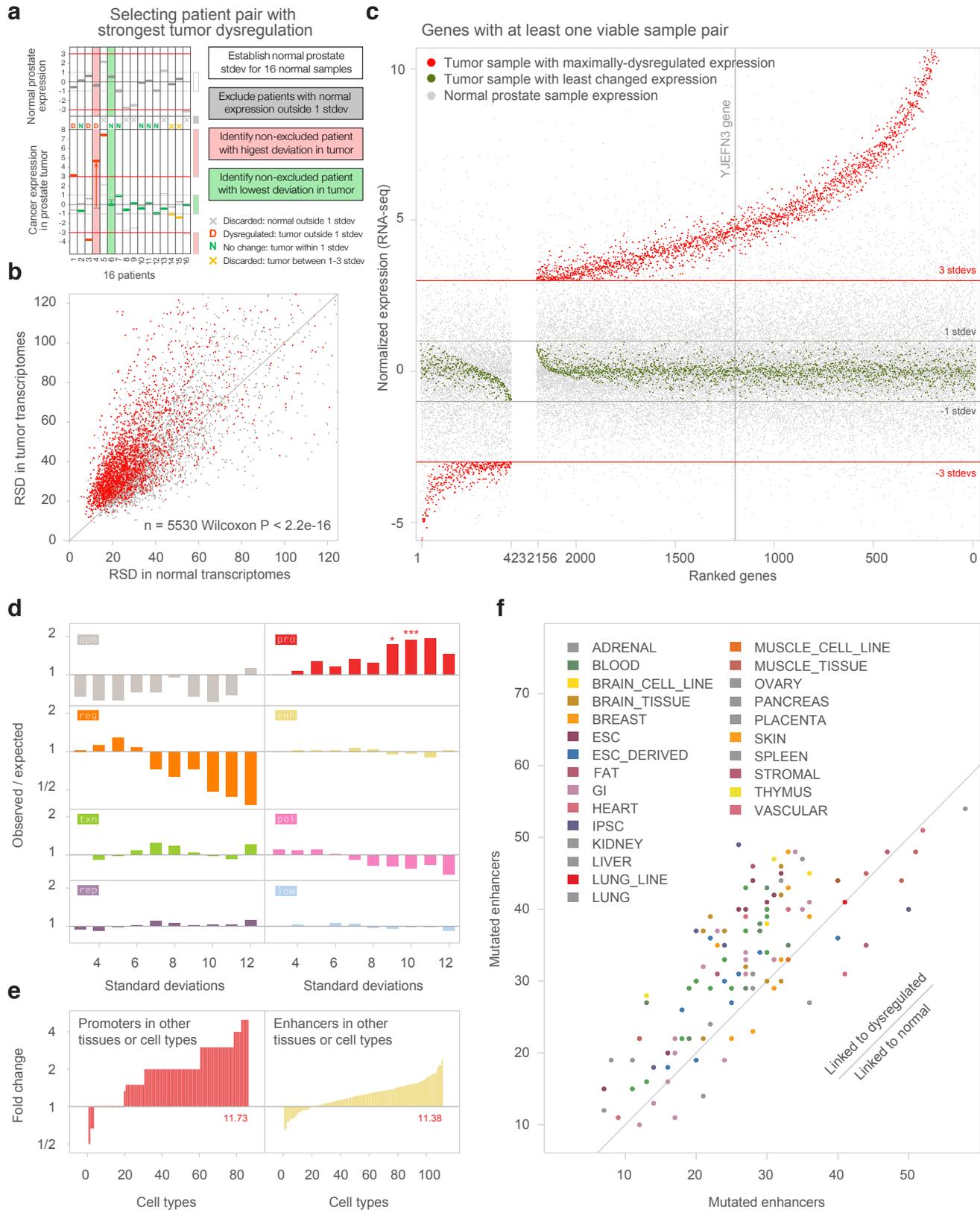


Figure 3

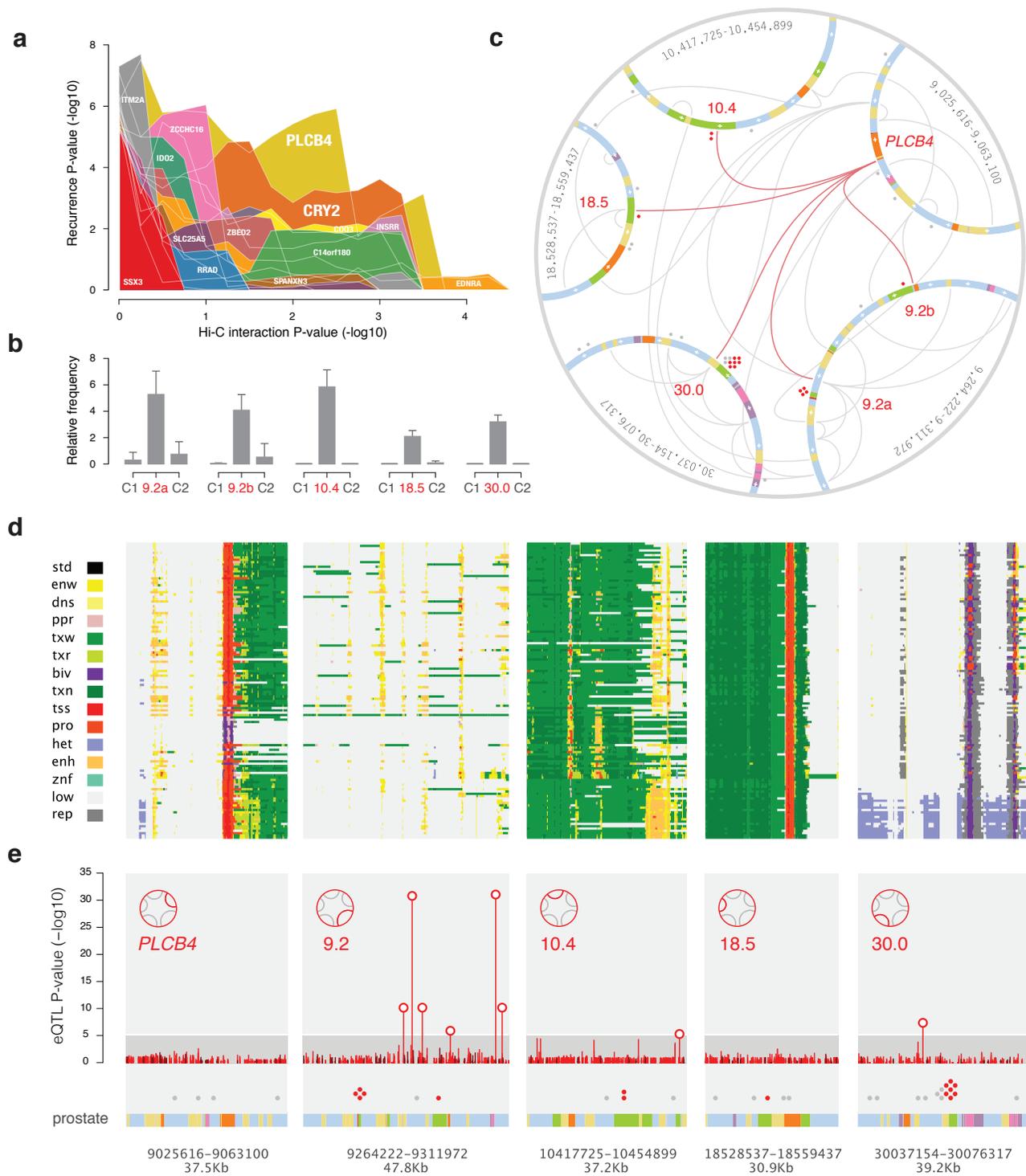


Figure 5

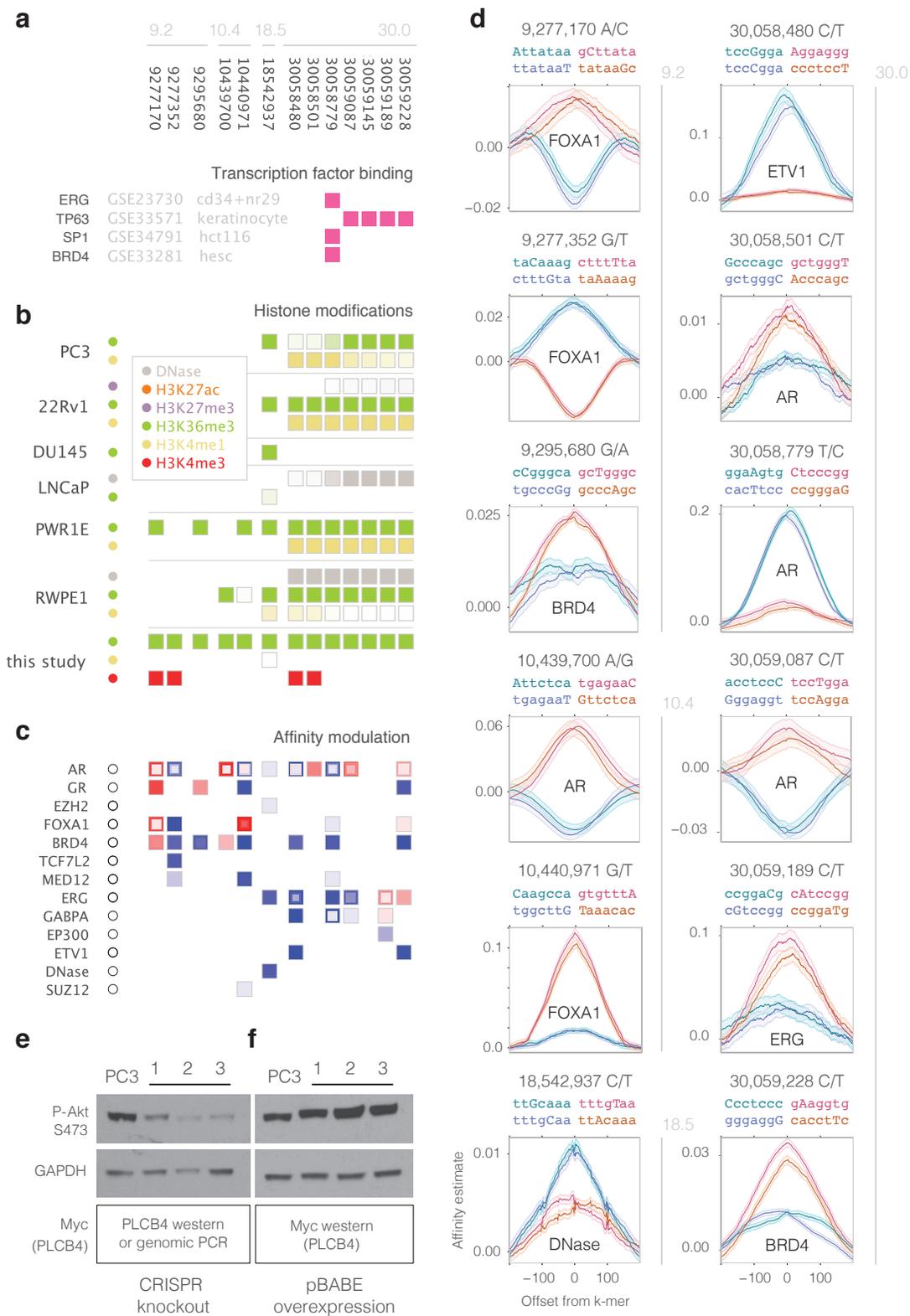
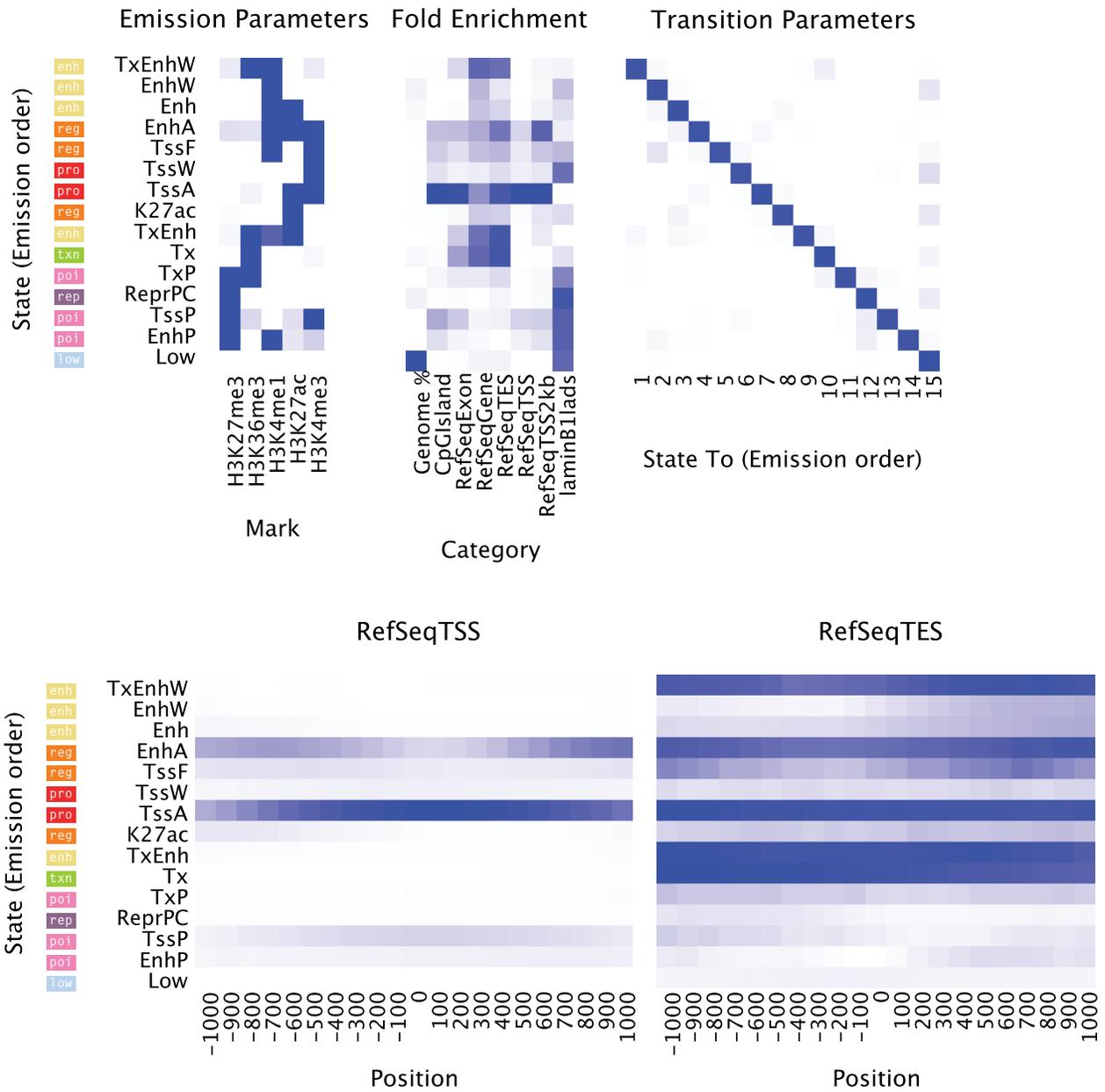
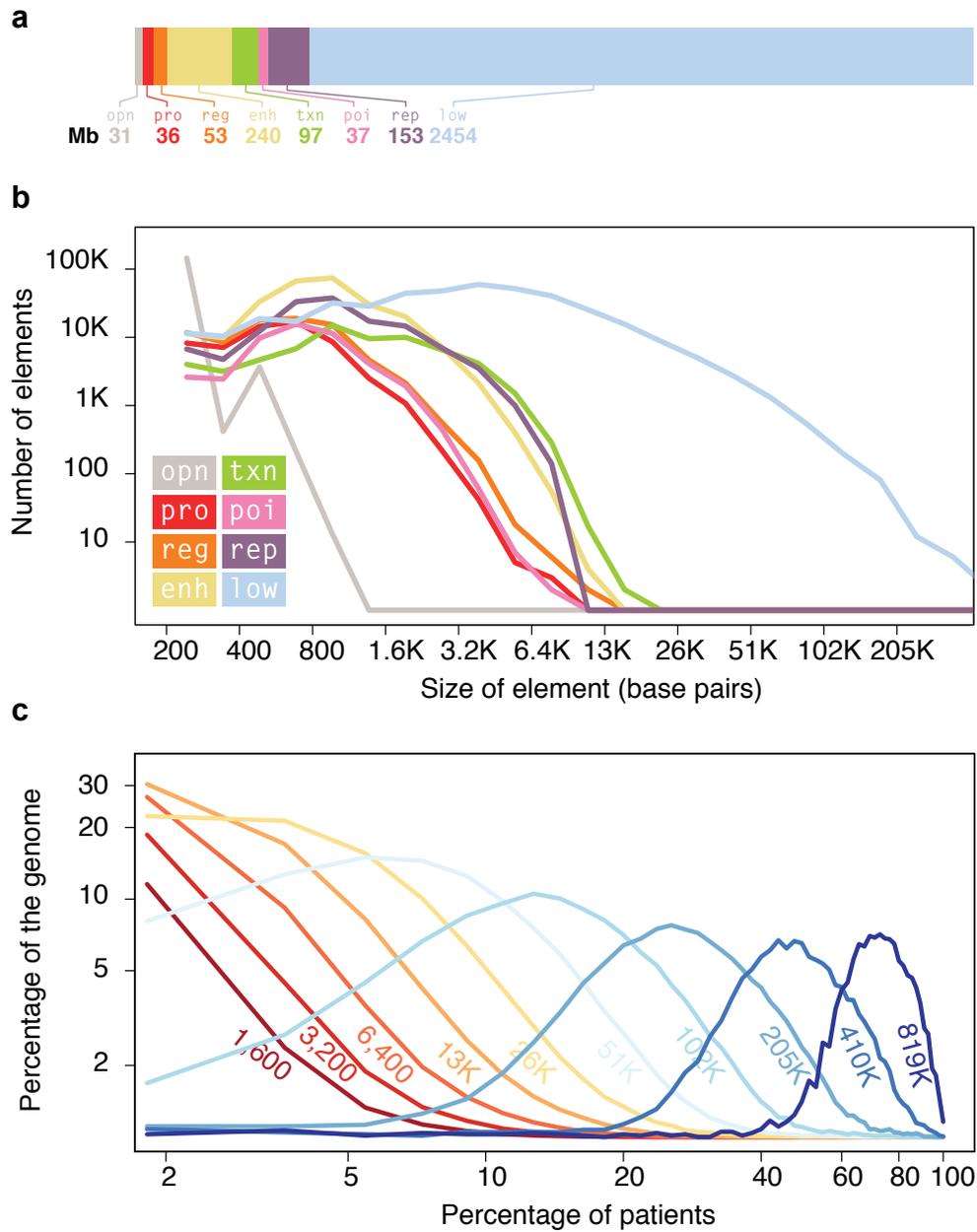


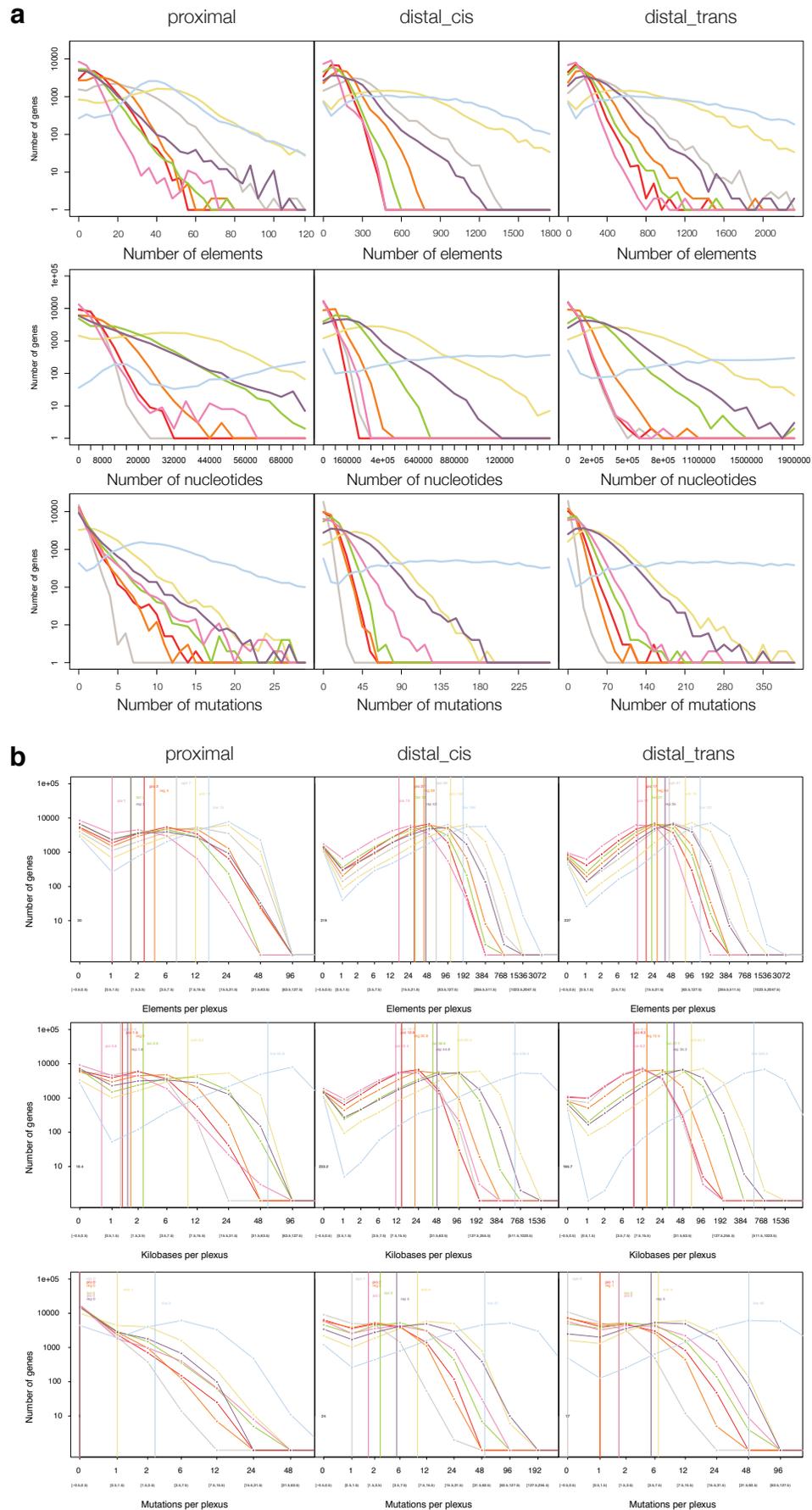
Figure 6



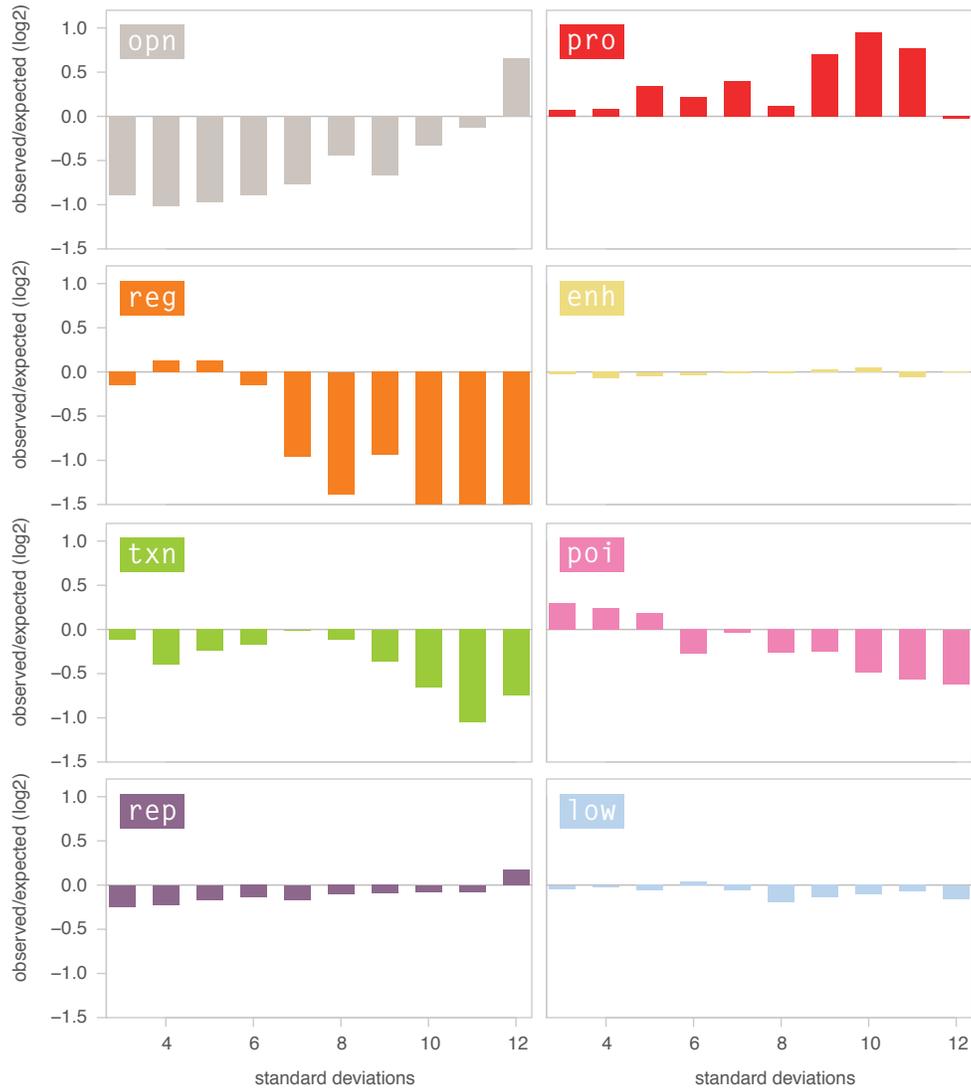
Supplementary Figure S1



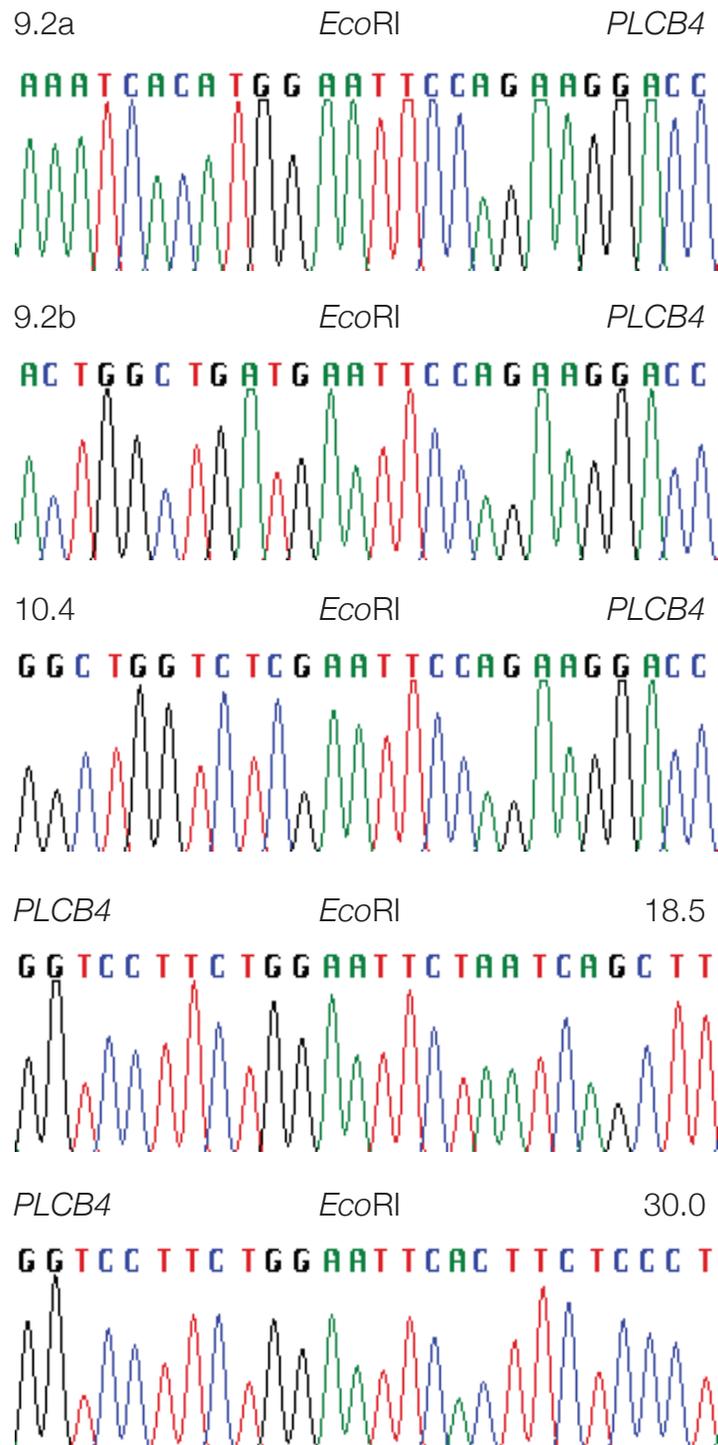
Supplementary Figure S2



Supplementary Figure S3



Supplementary Figure S4



Supplementary Figure S5

	opn	pro	reg	enh	txn	poi	rep	low
proximal	2.0	1.9	2.0	1.7	1.8	1.3	1.1	0.8
distal_cis	1.4	1.2	1.4	1.3	1.4	1.1	1.1	0.9
distal_trans	1.2	1.1	1.2	1.2	1.2	1.1	1.1	1.0

Supplementary Table 1 - Plexus of protein-coding genes are enriched for regulatory chromatin states

state	AA	AC	AG	AT	CA	CC	CG	CT	GA	GC	GG	GT	TA	TC	TG	TT	ALL	MutRteRnk	#MutRnk	CpG
opn	2038057	1524985	2309776	1556218	2265764	2443163	958272	2306343	1894140	2113434	2447267	1527792	1231077	1891959	2267317	2032368	30808000	7	8	958272
po	2850215	1813475	2504225	2088332	2557736	2697172	1134351	2497855	2178040	2208747	2704144	1820505	1670258	2167723	2568709	2839258	36302800	2	6	1134351
reg	4510920	2084157	3888669	3203227	3835214	3515950	1025193	3881141	3201168	2870862	3521743	2183325	2744768	3195527	3841497	4526761	53163100	5	7	1025193
enh	20782788	12158882	17706751	14843783	17795488	15074062	3605952	17703523	14331455	12624896	15115886	12198874	12582471	14321193	17842511	20830775	239524700	4	2	3605952
bn	8427032	5025053	7125216	6412020	7345157	5890525	1216510	7142206	5865845	4795317	5879282	503817	5351281	5883509	7353251	8479075	97228000	6	5	1216510
pol	2864388	1815799	2735960	2026880	2729322	2670900	864113	2737504	2310818	2207247	2676597	1827499	1538498	2307901	2745484	2867079	36927400	1	4	864113
rep	12465058	7717502	11671769	9430412	11807671	9982999	1872983	11672542	9687046	7945696	9987565	7738311	7324961	9690006	11826294	12496591	153321100	3	3	1872983
state	AA	AC	AG	AT	CA	CC	CG	CT	GA	GC	GG	GT	TA	TC	TG	TT	ALL	MutRteRnk	#MutRnk	CpG
opn	0.066	0.049	0.075	0.051	0.074	0.079	0.031	0.075	0.061	0.069	0.079	0.050	0.040	0.061	0.074	0.066	1.000	6	8	3.1%
po	0.079	0.050	0.069	0.058	0.070	0.074	0.031	0.069	0.060	0.061	0.074	0.050	0.046	0.060	0.071	0.078	1.000	2	6	3.1%
reg	0.085	0.050	0.073	0.060	0.072	0.066	0.019	0.073	0.060	0.054	0.066	0.050	0.052	0.060	0.072	0.085	1.000	4	7	1.9%
enh	0.087	0.051	0.074	0.062	0.074	0.063	0.015	0.074	0.060	0.053	0.063	0.051	0.053	0.060	0.074	0.087	1.000	4	2	1.5%
bn	0.087	0.052	0.073	0.066	0.076	0.061	0.013	0.073	0.060	0.049	0.060	0.052	0.055	0.061	0.076	0.087	1.000	5	5	1.3%
pol	0.078	0.049	0.074	0.055	0.074	0.072	0.023	0.074	0.063	0.060	0.072	0.049	0.042	0.062	0.074	0.078	1.000	1	4	2.3%
rep	0.081	0.050	0.076	0.062	0.077	0.065	0.012	0.076	0.063	0.052	0.065	0.050	0.048	0.063	0.077	0.082	1.000	3	3	1.2%
state	A	C	G	T	ALL	MutRteRnk	#MutRnk													
pol	9443030	9001040	9022162	9458966	36927400	1	4													
po	9256251	8887120	8911439	9245952	36302800	2	6													
reg	41284745	35336209	35358625	41337860	153321100	3	3													
enh	65492215	54179045	54271121	65576964	239524700	4	2													
reg	14319074	12260501	12277107	14302561	53163100	5	7													
bn	26989323	21594408	21574265	27067122	97228000	6	5													
opn	7428038	7973543	7982634	7422722	30808000	7	8													
state	A	C	G	T	ALL	MutRteRnk	#MutRnk	GC%												
bn	28%	22%	22%	28%	1	6	5	44%												
enh	27%	23%	23%	27%	1	4	2	45%												
reg	27%	23%	23%	27%	1	3	3	46%												
pol	27%	23%	23%	27%	1	5	7	40%												
po	26%	24%	24%	26%	1	1	4	49%												
po	25%	24%	25%	25%	1	2	6	49%												
opn	24%	26%	26%	24%	1	7	8	52%												

Supplementary Table 2 - Nucleotide and dinucleotide composition of chromatin states in prostate tissue

	Enriched chromatin state									All regulatory states								Size (kb)	
	Enrichment		Mutations			Fraction mutated		Convergence		Mutations		Fraction mutated		Convergence					
	state	P-val.	Nmut	N=1	N>1	Mut elmts	MutChrm	SinglElmt	Plexus	Nmut	N=1	N>1	Mut elmts	MutChrm	SinglElmt	Plexus	State	All	
ITM2A	enh	1E-07	58	9	7	16 / 245	9 / 20	13 24%	31 56%	129	18	18	36 / 793	12 / 20	14 25%	48 87%	181	483.4	
INSRR	poi	2E-07	76	8	9	17 / 142	11 / 18	12 22%	34 62%	267	49	40	89 / 2727	22 / 23	21 38%	55 100%	99.4	1756	
ZCCHC16	txn	3E-07	63	20	5	25 / 193	17 / 22	16 29%	33 60%	275	60	55	115 / 2138	23 / 23	18 33%	54 98%	251	1347	
ZBED2	pro	9E-07	72	23	12	35 / 331	13 / 21	10 18%	39 71%	421	74	69	143 / 4173	22 / 23	34 62%	55 100%	176	2771	
SPANXN3	pro	1E-06	24	2	5	7 / 76	4 / 14	5 9%	19 35%	124	31	24	55 / 1047	17 / 22	12 22%	47 85%	37.2	709.9	
PLCB4	txn	2E-06	79	9	15	24 / 341	14 / 23	10 18%	37 67%	362	59	54	113 / 3576	23 / 23	34 62%	54 98%	419	2497	
COQ3	pro	2E-06	58	19	9	28 / 247	11 / 20	9 16%	35 64%	326	57	58	115 / 2913	22 / 23	27 49%	55 100%	125	1957	
EDNRA	txn	2E-06	102	19	18	37 / 452	14 / 19	10 18%	41 75%	392	80	63	143 / 4238	21 / 23	42 76%	55 100%	586	2734	
CRY2	txn	3E-06	83	18	16	34 / 358	13 / 21	11 20%	36 65%	289	70	48	118 / 3569	20 / 23	33 60%	55 100%	455	2376	
ZC3H12B	rep	3E-06	150	38	24	62 / 576	16 / 22	20 36%	49 89%	415	75	66	141 / 2661	23 / 23	39 71%	55 100%	505	1754	
C14orf180	txn	4E-06	49	5	11	16 / 133	11 / 19	6 11%	28 51%	142	26	26	52 / 1700	20 / 23	30 55%	51 93%	161	1243	
IDO2	rep	4E-06	98	27	18	45 / 378	17 / 20	12 22%	48 87%	189	38	34	72 / 1525	19 / 21	26 47%	54 98%	368	1058	
RRAD	poi	4E-06	47	5	8	13 / 113	8 / 20	11 20%	28 51%	180	30	27	57 / 2091	19 / 22	30 55%	52 95%	71.8	1430	
SLC25A5	rep	4E-06	68	16	11	27 / 196	12 / 18	10 18%	37 67%	143	26	24	50 / 1332	17 / 21	16 29%	50 91%	189	852.7	
SSX3	enh	4E-06	53	6	8	14 / 230	9 / 20	14 25%	36 65%	104	14	17	31 / 682	17 / 21	14 25%	43 78%	196	468.4	

N=1 number of loci mutated in exactly one tumor

N>1 number of loci mutated in two or more tumors

Mut elmts number of mutated elements / Total number of elements in plexus

MutChrm number of mutated chromosomes / Total number of chromosomes in plexus

Size Total size of regulatory plexus for enrichment chromatin state, and for all regulatory states

Supplementary Table 3 - Plexus recurrence test results for all GENCODE protein-coding genes

internal_gene_id	internal_gene_id	internal_gene_id																	
G55697	ITM2A	enh	31	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
G03424	INSRR	poi	0	46	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
G56086	ZCCHC16	txn	0	0	41	0	0	0	0	0	0	0	0	0	0	0	0	0	0
G37166	ZBED2	pro	0	0	0	56	0	0	0	0	0	0	0	0	0	0	0	0	0
G56583	SPANXN3	pro	0	0	0	0	14	0	0	0	0	0	0	0	0	0	0	0	0
G32563	PLCB4	txn	0	0	0	0	0	54	5	0	0	0	0	0	0	0	0	0	0
G40609	EDNRA	txn	0	0	0	0	0	5	77	0	0	0	0	0	0	0	0	0	0
G45904	CQO3	pro	0	0	0	0	0	0	0	44	0	0	0	0	0	0	0	0	0
G08510	CRY2	txn	0	0	0	0	0	0	0	60	0	0	0	0	0	0	0	0	0
G20773	RRAD	poi	0	0	0	0	0	0	0	0	26	23	0	0	0	0	0	0	0
G20774	FAM96B(joined with RRAD)	poi	0	0	0	0	0	0	0	0	23	37	0	0	0	0	0	0	0
G55439	ZC3H12B	rep	0	0	0	0	0	0	0	0	0	0	125	0	0	0	0	0	0
G16752	C14orf180	txn	0	0	0	0	0	0	0	0	0	0	0	39	0	0	0	0	0
G50554	IDO2	rep	0	0	0	0	0	0	0	0	0	0	0	0	84	0	0	0	0
G56193	SLC25A5	rep	0	0	0	0	0	0	0	0	0	0	0	0	0	54	0	0	0
G55121	SSX3	enh	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	37	0

Supplementary Table 4 - Plexus intersections of all novel cancer-associated genes

Plexus loci	3C <i>Eco</i> RI primers	<i>Eco</i> RI fragment (chr20; hg19)	Sequence (5' to 3')
<i>PLCB4</i> promoter	Bait	9,044,304	gccattatggccttctcttggttcaactgtgg
Locus 9.2a	Control 1	9,266,126	ccaggagaactctttactgggtcacacagagtagc
	Test	9,272,249	cgcccaaatgttccacctgaagtcc
	Control 2	9,279,603	gagccttcacacatcccagttatgactgatcc
Locus 9.2b	Control 1	9,291,851	ggtgtggtgatgatgcaattctcctctgc
	Test	9,292,762	gatacaagaaagtcccatgggcaagaagaagg
	Control 2	9,296,754	cccattctcatttccctaagaaatgtcttggg
Locus 10.4	Control 1	10,429,887	gcaactagcaccatctgtcacctgtagaattgc
	Test	10,433,106	cgattctcctgtctcagcttcccgagtagc
	Control 2	10,440,835	ggaagtggaggatacaaaacttctgtttgaaagagc
Locus 18.5	Control 1	18,529,891	cgtatcgggtgggtccacatctataaattcaacc
	Test	18,541,239	ctggttgatcctcagcttaacaggcactgg
	Control 2	18,544,452	gcttttcattttctcaactgtgaccttcaaagagc
Locus 30	Control 1	30,049,488	ctggcttattcggattcttattgcacatatttgc
	Test	30,053,453	ccagtgttagtgaatctgtggcctaacttgtgacc
	Control 2	30,061,463	gcatcattcacttagactatgacatgcacgatgc

Supplementary Table 5 - *PLCB4* oligonucleotides