

# Pan-Cancer modelling of genomic alterations through gene expression

Federico M. Giorgi<sup>1, \*</sup> & Forest Ray<sup>2</sup>

<sup>1</sup> Department of Pharmacy and Biotechnology, University of Bologna, Via Selmi 3, 40138, Bologna, Italy

<sup>2</sup> Columbia University Medical Center, New York, NY 10032, USA

\* Corresponding Author

E-mail: federico.giorgi@unibo.it

# Short Title

## Pan-Cancer modelling of genomic alterations through gene expression

### Abstract

Cancer is a disease often characterized by the presence of multiple genomic alterations, which trigger altered transcriptional patterns and gene expression, which in turn sustain the processes of tumorigenesis, tumor progression and tumor maintenance. The links between genomic alterations and gene expression profiles can be utilized as the basis to build specific molecular tumorigenic relationships. In this study we perform pan-cancer predictions of the presence of single somatic mutations and copy number variations using machine learning approaches on gene expression profiles. We show that gene expression can be used to predict genomic alterations in every tumor type, where some alterations are more predictable than others. We propose gene aggregation as a tool to improve the accuracy of alteration prediction models from gene expression profiles. Ultimately, we show how this principle can be beneficial in intrinsically noisy datasets, such as those based on single cell sequencing.

### Author Summary

In this article we show that transcript abundance can be used to predict the presence or absence of the majority of genomic alterations present in human cancer. We also show how these predictions can be improved by aggregating genes into small networks to counteract the effects of transcript measurement noise.

# Introduction

Cancer is a molecular disease occurring when a cell or group of cells acquire uncontrolled proliferative behavior, conferred by a multitude of deregulations in specific pathways [1]. As is implied by such a broad definition, cancer is a highly heterogeneous disease, showing remarkably different molecular, histological, genetic and clinical properties, even when comparing tumors originating from the same tissue [2]. Many cancers are characterized by the presence of single nucleotide or short indel mutations and/or copy number alterations, which appear somatically at the early stages of oncogenesis and can drive tumor progression [3]. Cancers can be broadly divided in two classes: the M class, where point mutations are prevalent, and the C class, where copy number variations (CNVs) are more numerous and are often associated with TP53 mutations. Tumor class influences anatomic location. Most ovarian cancers, for example, belong to the C class, while most colorectal cancers belong to the M class, although many exceptions do exist [4].

The Cancer Genome Atlas (TCGA) project [5] has recently underwent a major effort to collect vast amounts of information on thousands of distinct tumor samples. The TCGA data collection, commonly referred to as the “Pan-cancer” dataset, provided the scientific community with an avalanche of data on DNA alterations, gene expression, methylation status and protein abundances among others, with the critical mass necessary to identify rarer driver tumorigenesis effects in many types of cancers [6–8]. By combining all 33 TCGA datasets, Bailey and colleagues [9] recently outlined a pan-cancer map of which mutations can be drivers for the progression of cancer.

The availability of thousands of samples measuring many different variables in cancer has allowed scientists to generate statistical models of relationships between different

molecular species. A pan-cancer correlation network between coding genes and long noncoding RNAs, for example, sheds light on the function of non-coding parts of the transcriptome [10]. More recently, mutations on transcription factors (TFs) have been linked to altered gene expressions and phosphoprotein levels in 12 TCGA tumor type datasets [11]. Network approaches have been applied to identify clusters of coexpressed genes, shared by multiple cancer types [12]. Several studies have sought to characterize the relationships between genomic status and expression levels in cancer, trying to identify commonalities across different cancer types [13,14]. In particular, Alvarez and colleagues [15] have postulated that the effect of genomic alterations in cancer can be more readily assessed by aggregating gene expression profiles into transcriptional networks, rather than by profiles taken separately.

While the association between genomic events and gene expression is proven in several scenarios, it remains to be seen if it can be assessed in scenarios where fully quantitative readouts are unavailable, such as low coverage samples. One of these scenarios is Single Cell Sequencing [16], often carried out in experiments where thousands of mutations are generated via a system of pooled CRISPR-Cas9 knockouts [17].

To our knowledge there is no study trying to identify relationships between all genomic alteration events (somatic mutations/indels and CNVs) and global gene expression across cancers. In this study, we use 24 TCGA tumor datasets to investigate whether gene expression can be used to predict the presence of specific genomic alterations in several cancer tissue contexts. To this end, we leverage the current availability of a vast family of machine learning algorithms [18]. We investigate whether some gene alterations can be better modelled than others, and whether using grouped gene expression profiles as

aggregated variables can effectively identify specific genomic alterations. Finally, we test whether predicting mutations and CNVs can be carried out in an intrinsically noisy single cell RNA-Seq (scRNA-Seq) transcriptomics datasets.

## Results

### Collection of Pan-Cancer Dataset

We downloaded the most recent version of the TCGA datasets available on Firehose (v2016\_01\_28), encompassing mutational, CNV and gene expression data. Using TSNE clustering on gene expression data (9642 samples), we observed how different tumor types cluster separately from each other (Figure 1A). However, two tumour types segregate into two subgroups: breast cancer, which subdivides into a major luminal cluster and a smaller (in terms of samples collected) basal cluster [19]; and esophageal carcinoma, which roughly subdivides into adenocarcinomas and squamous cell carcinomas [20].

We then aggregated the single nucleotide and short indel somatic mutation data from the same samples for which we had collected gene expression. As is widely known, TP53 is the most mutated gene in human cancer (Figure 1B), followed by PIK3CA, SYNE1 and KRAS. As shown before [4] some tumor types are characterized by a high presence of somatic mutations. In particular, colorectal cancer, mesothelioma and esophageal cancer carry at least one of these events in almost 100% of the samples in the TCGA dataset. In the figure, we filtered out commonly known non-driver mutations [21], such as those happening in long genes like TTN and OBSCN, but we kept them in all following analyses for the sake of completion. A representation of all mutated genes, including blacklisted ones, is available in

Figure S1. Some tumors are characterized by the prevalence of a mutation in a specific gene, such as the G-protein coding BRAF in thyroid carcinoma [22] or IDH1, translating into isocitrate dehydrogenase, in low grade glioma [23].

Finally, we obtained readouts of CNV status for all TCGA samples. CNVs can have different extensions in terms of nucleotides affected and can sometimes encompass entire chromosomes [24] and the thousands of genes therein. In order to limit the number of variables to a more meaningful subset, we assigned a CNV profile to every gene, and kept only those whose CNV profiles are positively and significantly correlated with their transcript abundance profiles [25]. We defined these events as functional CNVs (fCNVs). In order to make fCNV variables comparable to the mutational ones, we defined a cut-off for presence or absence by using the  $\log_2(\text{CNV})$  threshold of 0.5, which roughly corresponds to at least one copy gain for amplifications, and at least one copy loss for deletions (see Materials and Methods). We then reported their abundance in the pan-cancer dataset, distinguishing between amplifications (Figure 1C) and deletions (Figure 1D). As previously shown [4], virtually all ovarian cancer samples are characterized by at least one CNV event. Among the most amplified genes, we find the oncogenes SOX2 [26], EGFR [27] and MDM2 [28], and also a non-coding gene, PVT1, the most amplified gene in breast cancer, with proven but as-of-yet uncharacterized proto-oncogenic effects [29,30]. Amongst the most deleted genes (Fig.1D) we observe well known tumor-suppressor genes, such as CDKN2A [31,32] and PTEN [33,34].

1

## 2   Modelling Cancer Alterations with gene expression

3   After collecting all the expression and genomic alteration data from TCGA, we set out to  
4   generate models able to predict the presence or absence of each event by virtue of gene  
5   expression data in the contexts of all collected tumor types.

6   We tested several modelling algorithms for classification using the aggregator platform for  
7   machine learning caret [18] in the bladder cancer mutational dataset [35]. We observed that  
8   all models provide better-than-random predictions for the majority of mutational events, in  
9   terms of area under the ROC curve (AUROC)(Figure 2) [36]. We chose the top-scoring  
10   algorithm in this test, the Gradient Boost Modelling algorithm (gbm), a robust tree-based  
11   boosting model [37], due to its robustness and speed of implementation.

12   We calculated gbm models for all tumour types of at least 100 samples with co-measured  
13   expression and CNV or mutations, which included 24 of the 33 TCGA tumor types. The  
14   models were predictive of genomic events observed in no less than 5% and no more than  
15   95% of the patients in the dataset, and at least in 10 samples. Our results show that in all  
16   tumour types, a machine learning algorithm based on gene expression is consistently better  
17   than a random predictor (AUROC line at 0.5) at correctly classifying tumour samples for the  
18   presence or absence of specific genomic alteration events (Figure 3 and Supplementary  
19   Table S1). In particular, TP53 mutations are well modelled in many of these tumor types,  
20   being the most well predicted mutational event in both acute myeloid leukemia and low  
21   grade glioma. We could also model the presence of a copy loss of TP53 in sarcoma, which  
22   can be predicted with an accuracy of 70%. Ovarian and pancreatic cancer datasets  
23   presented exceptional cases, in that each contained such high TP53 mutation rates (next to

95% detected) [38,39] that our algorithms could not distinguish sufficient differences within each dataset to train a model. Also KRAS-targeting events are well modelled, specifically in colon, lung and stomach cancer, and cervical squamous carcinoma [40]. We noted a tendency where models for more frequent CNV events yielded a greater predictive power (Figure S2), a tendency not observed for somatic mutation models. We then tested if known tumor-related genes, such as those curated by the Cancer Gene Census [41] are better modelled than the rest of the genome. There is no difference in mutation and amplification results, but for deletion events, oncogenes yield weaker models (Wilcoxon Test,  $p=0.0037$ ) and tumor suppressor genes yield generally stronger models ( $p=0.00050$ ). This is in agreement with the central paradigm of cancer, where a tumor suppressor gene deletion can be one of the driving events of tumorigenesis and tumor progression [42]. On the other hand, deletion of tumour-promoting oncogenes is generally unfavourable for tumor progression, and so, generally speaking it should be present only as a passenger event, unlikely to determine global gene expression and tumor fate.

## Modelling specific alterations with noise addition

In order to understand whether cancer-related genomic alterations can be modelled by gene expression in scenarios with lower signal-to-noise ratio, we artificially perturbed the TCGA gene expression dataset via the addition of Gaussian noise, and then proceeded to build models to predict the presence of TP53 mutations in breast cancer, the largest dataset in TCGA by number of samples.



As expected, the addition of uniform random gaussian noise to the gene expression matrix has a detrimental effect on the amount of information left for modelling the presence of TP53 somatic mutations (Figure 4A).

We then decided to test several permutations of noise addition on the same breast cancer expression data, by each time aggregating genes into networks defined a priori in the same context, using a Tukey Biweight Robust Average method [43] on Weighted Gene Correlation Network Analysis (WGCNA) clusters [44] and the VIPER algorithm [15] on ARACNe-AP networks [45]. It is important to note that WGCNA clusters are completely non-overlapping and yield generally a lower number of aggregated variables than VIPER clusters, which are groups of genes possibly shared by other transcription factor clusters and that collectively yield the global expression of a transcription factor target set (dubbed as a proxy for “TF activity” in the original VIPER manuscript [15]).

Our results show that gene expression, VIPER activity and WGCNA clusters yield very similar models for predicting TP53 mutations in breast cancer (figure S4). The amount of information contained in the input variables is therefore comparable. Adding noise to the input expression matrix, however, and then aggregating the resulting noise-burdened genes into VIPER or WGCNA clusters (see Materials and Methods), provides robustness to the models (Figure 4B). Similar results with higher variances (possibly due to the smaller size of the datasets) can be observed for EGFR amplifications in glioblastoma (Figure S5) and lung squamous carcinoma (Figure S6), for PVT1 amplifications in ovarian cancer (Figure S7) and for PTEN deletions in sarcoma (Figure S8). In all these examples, however, the performance of the simple WGCNA/Tukey aggregation is closer (if not worse) to that of simple gene expression.

An alternative way to reduce the information content from an NGS gene expression dataset is to reduce the number of read counts from each sample. This operation reflects either a low coverage bulk RNA-Seq experiment or an experiment arising from Single-Cell sequencing [46]. In particular, single-cell RNA-Seq (scRNA-Seq) is characterized by the dropout phenomenon [47] wherein genes expressed in the cells are sometimes not detected at all. In order to simulate such scenarios, we down-sampled each RNA-Seq gene count profile from the largest TCGA dataset (Breast Cancer) to a target aligned read number using a beta function, which allows for reduction coupled with random complete gene dropouts (Figure 5A). We then modelled again the presence of TP53 mutations using gene expression (Figure 5B). We found out that models based on standard unaggregated gene expression experience an accuracy drop at around 30M reads, while aggregating genes using VIPER (but not with WGCNA) allows for better-than-random accuracies even at 3M reads, confirming the benefits of gene aggregation in low coverage RNA-Seq, as previously found e.g. for sample clustering [48].

## Mutation prediction in single-cell data

We set out to detect if mutations can be modelled from gene expression data in single-cell RNA-Seq contexts. In order to do so, we used the original CROP-Seq dataset [17], where multiple gene knock-outs were carried out via CRISPR/Cas9 in Jurkat cells and the presence of the deletion was measured alongside gene expression in a single cell manner.

We built models based on 8 knock-out subsets targeting the following genes: JUNB, JUND, LAT, NFAT5, NFKB1, NFKB2, NR4A1 and PTPN11, all with at least 35 single cells carrying the single knock-outs (vs. 420 control wild-type single cells). Our analysis shows that gene

aggregation in TF-centered coexpression groups using ARACNe/VIPER can be beneficial in predicting mutation presence, by virtue of showing the probability of carrying the mutation in mutated samples vs. control samples (Figure 6).

## Discussion

In this paper, we tested a framework to investigate the complex relationships between genetic events and transcriptional deregulation through machine learning approaches. We demonstrated as a generalized proof-of-principle that genomic alterations can be modeled by gene expression across several human cancers through several machine learning algorithms, and specifically that a gradient boost modeling approach seems optimal for the task. In the process, we generated a collection of models for each genomic alteration in each cancer context, showing that the best predicted alterations are not necessarily targeting known oncogenes or tumor suppressors. Interestingly, we show how the aggregation of gene expression profiles in groups of coexpressed genes, via the ARACNe/VIPER or WGCNA methods, makes the models more robust and more resistant to perturbations such as gaussian noise or artificial downsampling. Finally, we have shown how the same aggregation principle can have beneficial effects in predicting the presence of mutations in intrinsically noisy scenarios, like single cell RNA-Seq. At the same time, we have shown how modeling can be carried out in single-alteration contexts, implicitly overtaking the potential bias of cancer samples, where in fact multiple genomic alterations can and do coexist.

The performance of gene aggregation methods has been tested before for sample clustering in RNA-Seq read reduction scenarios [15,48], but never in this specific task nor in a pan-cancer context. As a principle, the usage of robust averages of pre-defined co-expressed

genes can be applied in any context where reliability of gene expression data is necessary, from differential expression to pathway enrichment analyses. The notion that relationships between genomic alterations and gene expression profiles can be robustly modelled across different cancer scenarios, as well as in single-cell and noisy contexts, can have important repercussions in diagnostics, where theoretically a single quantitative expression experiment can be used to predict the presence or absence of a mutation.

## Materials and Methods

### Data processing

We obtained raw expression counts, mutation and CNV raw data from TCGA using the Firehose portal ([gdac.broadinstitute.org](http://gdac.broadinstitute.org)). Raw counts were normalized using Variance Stabilizing Transformation as described before [49]. Somatic mutations not changing the aminoacid sequence of the protein product were discarded. We flagged genes blacklisted by the MutSig project [21], such as TTN, ORs, MUCs as false positives, and removed them from further analysis (except the most mutated in the pan-cancer dataset, shown in Figure S1). CNV tracks were associated to the targeted gene using the GenomicRanges R package [50]. Gene-centered CNVs were then associated to the expression profile of the gene itself. CNV tracks with a Spearman correlation coefficient above 0.5 were deemed “functional CNVs” [25] and used in the rest of the analysis. Samples with more than 0.5% of the genes in the genome somatically amplified, deleted or mutated were deemed “hypermodified” and the total number was shown in Figure 1 bottom bars.

Clustering analysis was carried out on the TCGA tumor samples using the expression profiles of 1172 Transcription Factors defined by Gene Ontology terms “transcription factor activity, sequence-specific DNA binding” (GO:0003700) and “nuclear location” (GO:0005634) [51].

The dataset expression profiles were visualized after TSNE transformation [52] with 1000 iterations using a 2D kernel density estimate for coloring different tumor types [53].

Oncogenes and Tumor Suppressor genes were obtained from the COSMIC Cancer Gene Census in October 2018 [41].

## Modeling

We used the R caret package [18] as the platform to run all our predictive models in a standardized and reproducible way. Binary classifiers were built to predict the presence/absence of mutation, amplification and deletion events. The CNV value provided by TCGA corresponds to  $\log_2(\text{tumor coverage}) - \text{genomic median coverage}$ . The threshold for amplification/deletion presence was set to 0.5.

Data partitioning was performed once for each tumor type, with 75% of the samples used for training and 25% for test purposes. Training was performed using 10-fold Cross Validation. Recursive Feature Elimination was carried out by the default caret implementation on the 10,000 highest variance gene expression tracks. The algorithms used (and R packages implementing theme) were:

- Bayesian Generalized Linear Model (bayesglm)
- Tree Models from Genetic Algorithms (evtree)
- Gradient Boost Modeling (gbm)
- Generalized Linear Model (glm)

- k-Nearest Neighbors (knn)
- Linear Discriminant Analysis (lda)
- Neural Networks (mxnet)
- Neural Networks with Feature Extraction (pcaNNet)
- Random Forest (rf)
- Linear Support Vector Machine (svmLinear)
- Radial Support Vector Machine (svmRadial)

In order to reduce information from the gene expression profiles, we adopted two strategies. The first, shown e.g. in Figure 4B, adds random gaussian noise to the expression tracks, with a variable standard deviation (indicated as “Gaussian Noise Level”). Each model run after noise addition was run 100 times to allow for various data partitions. The second strategy (Figure 5) reduced the number of reads mapped to each gene in order to obtain expression samples with decreased total gene counts. In order to do so, we applied to each gene in each sample a downsampling factor sample from a beta distribution:

$$\frac{1}{B(\alpha, \beta)} x^{\alpha-1} (1-x)^{\beta-1}$$

Where B is the Beta function, acting as a normalization constant, x is the raw gene expression count in a particular sample,  $\alpha$  is the first shape parameter and  $\beta$  the second shape parameter. In order to reduce the total sample coverage to the desired level,  $\beta$  is set to 0.1 and  $\alpha$  is set to:

$$\alpha = \frac{\beta f/r}{1 - f/r}$$

Where  $f$  is the desired number of reads and  $r$  is the total number of reads in the sample. A real case example of this beta distribution is shown in Figure S9.

### Aggregation algorithms

We used ARACNe-AP [45] to generate TF-centered networks on each of the VST-normalized TCGA expression datasets. TFs were selected via Gene Ontology as described before, with  $p$ -value for each network edge set to  $10^{-8}$ . ARACNe networks were then used to obtain an aggregated value of TF activity for each sample using the VIPER algorithm [15] which reports the collective gene expression level changes of each TF-centered network vs. the mean expression of each gene in the dataset. Only TF networks with at least 10 genes (excluding the TF) were included.

WGCNA clusters of genes were constructed using the *wcna* package [44] with default parameters and minimum network size set to 10. To obtain a robust median expression value for each WGCNA cluster in each sample we used Tukey's Biweight function as implemented by the R *affy* package [54].

### Single Cell dataset

CROP-Seq raw expression counts were obtained from the Datlinger dataset (available on Gene Expression Omnibus, entry GSE92872). Samples mapping wild-type control cells and the most represented knock-out genes (JUNB, JUND, LAT, NFAT5, NFKB1, NFKB2, NR4A1 and PTPN11) were selected. Variance Stabilizing Transformation was applied using a blinded experimental design. Gradient boost modelling was applied to each model as described in the previous paragraph, and probabilities of carrying the knock-out for samples in the test set are shown, grouped for wild-type and knock-out samples. In this particular case, 10 data

partitioning rounds are done, in order to increase the exploration space of the model performance.

## Methods Availability

All code used to generate the analysis and the figures of this paper is available in the online materials.

## Figure Legends

**Figure 1.** The TCGA dataset used. A: TSNE clustering of TCGA samples based on the expression profiles of Transcription Factors. The 2D median of each tumor type is indicated using the TCGA tumor code. Subset size is indicated in brackets next to tumor type names to the right. B: table of most somatically mutated genes across TCGA tumor samples, in terms of number of samples where the gene is somatically mutated with altered protein product sequence. C: table of most amplified genes across TCGA tumor samples. D: table of most deleted genes across TCGA tumor samples. The fraction of total TCGA samples carrying a gene-targeting event is indicated to the right of panels B-D, and the fraction of samples where more than 0.5% of the genes is affected by the panel event type is indicated to the bottom of panels B-D.

**Figure 2.** Performance of 11 machine learning algorithms in binary classification of mutated/nonmutated samples using gene expression predictor variables in the Bladder Cancer dataset. Each point corresponds to a specific mutation/model. Performance is indicated as AUROC: Area Under the Receiver Operating Characteristic curve.



**Figure 3.** Performance of gbm models for each genomic alteration event in TCGA, predicted as a function of each tumor gene expression. Alterations targeting TP53 and KRAS are indicated.

**Figure 4.** Performance of a TP53 somatic mutation gbm model upon gaussian noise addition. A: ROC curves (and AUC) upon addition of increasing levels (in terms of SD of a gaussian distribution with mean=0) of gaussian noise. B: AUROCs of the model with increasing noise, calculated using gene expression (black line) or aggregated gene expression using the WGCNA (green line) or VIPER (red line) algorithms. Pseudocounts of 0.1 are added in order to show zero counts as -1 in  $\log_{10}$  scale.

**Figure 5.** Performance of a TP53 mutation gbm model upon downsampling of the TCGA breast cancer RNA-Seq dataset. A: for a single TCGA sample (TCGA-A1-A0SB-01) with 43.8 gene mapping reads, the downsampling algorithm is applied for multiple target read quantities. X-axis shows the count for each gene in the original sample, and Y-axis in the downsampled output. B: AUROCs of the model with decreasing read numbers, calculated using gene expression (black line) or aggregated gene expression using the WGCNA (green line) or VIPER (red line) algorithms.

**Figure 6.** Modeling of single cell KO mutations using single cell gene expression in the Datlinger dataset. Each point indicates a sample in multiple test sets. Known Wild Type Control samples (CTRL, left) are plotted separately from Known Knock-Out samples (KO, right), with number in brackets indicating the number of cells carrying the specific genotype. The probability of carrying a mutation is shown on the y axis. Boxplots showing median distribution are overlaid on the sample KO probability distributions. Results using standard VST-normalized expression data are shown (green) for each gene next to identical models

run with aggregated gene expression using the VIPER algorithm (blue). One-tailed Wilcoxon tests were calculated between the KO and CTRL distributions of probabilities, and p-values are reported.

## Acknowledgments

We thank Dr. Marco Russo and Jordan Pflugh Kraft for the fruitful discussions.

## References

- [1] Hanahan D, Weinberg RA. Hallmarks of cancer: the next generation. *Cell* 2011;144:646–74. doi:10.1016/j.cell.2011.02.013.
- [2] Meacham CE, Morrison SJ. Tumour heterogeneity and cancer cell plasticity. *Nature* 2013;501:328–37. doi:10.1038/nature12624.
- [3] Bozic I, Antal T, Ohtsuki H, Carter H, Kim D, Chen S, et al. Accumulation of driver and passenger mutations during tumor progression. *Proc Natl Acad Sci U S A* 2010;107:18545–50. doi:10.1073/pnas.1010978107.
- [4] Ciriello G, Miller ML, Aksoy BA, Senbabaoglu Y, Schultz N, Sander C. Emerging landscape of oncogenic signatures across human cancers. *Nat Genet* 2013;45:1127–33. doi:10.1038/ng.2762.
- [5] Chang K, Creighton CJ, Davis C, Donehower L, Drummond J, Wheeler D, et al. The Cancer Genome Atlas Pan-Cancer analysis project. *Nature* 2013.
- [6] Brennan CW, Verhaak RGW, McKenna A, Campos B, Nounshmehr H, Salama SR, et al. The Somatic Genomic Landscape of Glioblastoma. *Cell* 2013;155:462–77. doi:10.1016/j.cell.2013.09.034.
- [7] Cancer Genome Atlas Network. Genomic Classification of Cutaneous Melanoma. *Cell* 2015;161:1681–96. doi:10.1016/j.cell.2015.05.044.
- [8] Leiserson MDM, Vandin F, Wu H-T, Dobson JR, Eldridge JV, Thomas JL, et al. Pan-cancer network analysis identifies combinations of rare somatic mutations across pathways and protein complexes. *Nat Genet* 2015;47:106–14. doi:10.1038/ng.3168.
- [9] Bailey MH, Tokheim C, Porta-Pardo E, Sengupta S, Bertrand D, Weerasinghe A, et al. Comprehensive Characterization of Cancer Driver Genes and Mutations. *Cell* 2018;174:1034–5. doi:10.1016/j.cell.2018.07.034.
- [10] Liu Y, Zhao M. InCaNet: pan-cancer co-expression network for human lncRNA and cancer genes. *Bioinformatics* 2016;32:1595–7. doi:10.1093/bioinformatics/btw017.
- [11] Osmanbeyoglu HU, Toska E, Chan C, Baselga J, Leslie CS. Pancancer modelling predicts the context-specific impact of somatic mutations on transcriptional programs. *Nat Commun* 2017;8:14249. doi:10.1038/ncomms14249.
- [12] Kim H, Kim Y-M. Pan-cancer analysis of somatic mutations and transcriptomes reveals common functional gene clusters shared by multiple cancer types. *Sci Rep* 2018;8:6041. doi:10.1038/s41598-018-24379-y.

- [13] Ghazanfar S, Yang JYH. Characterizing mutation-expression network relationships in multiple cancers. *Comput Biol Chem* 2016;63:73–82. doi:10.1016/j.compbiolchem.2016.02.009.
- [14] Sharma A, Jiang C, De S. Dissecting the sources of gene expression variation in a pan-cancer analysis identifies novel regulatory mutations. *Nucleic Acids Res* 2018;46:4370–81. doi:10.1093/nar/gky271.
- [15] Alvarez MJ, Shen Y, Giorgi FM, Lachmann A, Ding BB, Ye BH, et al. Functional characterization of somatic mutations in cancer using network-based inference of protein activity. *Nat Genet* 2016;48:838–47. doi:10.1038/ng.3593.
- [16] Nawy T. Single-cell sequencing. *Nat Methods* 2013;11:18. doi:10.1038/nmeth.2771.
- [17] Datlinger P, Rendeiro AF, Schmidl C, Krausgruber T, Traxler P, Klughammer J, et al. Pooled CRISPR screening with single-cell transcriptome readout. *Nat Methods* 2017;14:297–301. doi:10.1038/nmeth.4177.
- [18] Kuhn M. Building Predictive Models in R Using the caret Package. *J Stat Softw* 2008;028.
- [19] Perou CM, Sørlie T, Eisen MB, van de Rijn M, Jeffrey SS, Rees CA, et al. Molecular portraits of human breast tumours. *Nature* 2000;406:747–52. doi:10.1038/35021093.
- [20] Not All Esophageal Tumors Equal. *Cancer Discov* 2017;7:238. doi:10.1158/2159-8290.CD-NB2017-006.
- [21] Lawrence MS, Stojanov P, Polak P, Kryukov GV, Cibulskis K, Sivachenko A, et al. Mutational heterogeneity in cancer and the search for new cancer-associated genes. *Nature* 2013;499:214–8. doi:10.1038/nature12213.
- [22] Kimura ET, Nikiforova MN, Zhu Z, Knauf JA, Nikiforov YE, Fagin JA. High Prevalence of BRAF Mutations in Thyroid Cancer: Genetic Evidence for Constitutive Activation of the RET/PTC-RAS-BRAF Signaling Pathway in Papillary Thyroid Carcinoma. *Cancer Res* 2003;63:1454–7.
- [23] Yan H, Parsons DW, Jin G, McLendon R, Rasheed BA, Yuan W, et al. IDH1 and IDH2 Mutations in Gliomas. *N Engl J Med* 2009;360:765–73. doi:10.1056/NEJMoa0808710.
- [24] Shlien A, Malkin D. Copy number variations and cancer. *Genome Med* 2009;1:62. doi:10.1186/gm62.
- [25] Chen JC, Alvarez MJ, Talos F, Dhruv H, Rieckhof GE, Iyer A, et al. Identification of Causal Genetic Drivers of Human Disease through Systems-Level Analysis of Regulatory Networks. *Cell* 2014;159:402–14. doi:10.1016/j.cell.2014.09.021.
- [26] Bass AJ, Watanabe H, Mermel CH, Yu S, Perner S, Verhaak RG, et al. SOX2 is an amplified lineage-survival oncogene in lung and esophageal squamous cell carcinomas. *Nat Genet* 2009;41:1238–42. doi:10.1038/ng.465.
- [27] Bell DW, Lynch TJ, Haserlat SM, Harris PL, Okimoto RA, Brannigan BW, et al. Epidermal Growth Factor Receptor Mutations and Gene Amplification in Non-Small-Cell Lung Cancer: Molecular Analysis of the IDEAL/INTACT Gefitinib Trials. *J Clin Oncol* 2005;23:8081–92. doi:10.1200/JCO.2005.02.7078.
- [28] Momand J, Jung D, Wilczynski S, Niland J. The MDM2 gene amplification database. *Nucleic Acids Res* 1998;26:3453–9. doi:10.1093/nar/26.15.3453.
- [29] Colombo T, Farina L, Macino G, Paci P. PVT1: a rising star among oncogenic long noncoding RNAs. *BioMed Res Int* 2015;2015:304208. doi:10.1155/2015/304208.
- [30] Li X, Chen W, Wang H, Wei Q, Ding X, Li W. Amplification and the clinical significance of circulating cell-free DNA of PVT1 in breast cancer. *Oncol Rep* 2017;38:465–71. doi:10.3892/or.2017.5650.

- [31] Usvasalo A, Savola S, Rätty R, Vettenranta K, Harila-Saari A, Koistinen P, et al. CDKN2A deletions in acute lymphoblastic leukemia of adolescents and young adults—An array CGH study. *Leuk Res* 2008;32:1228–35. doi:10.1016/j.leukres.2008.01.014.
- [32] Mistry M, Zhukova N, Merico D, Rakopoulos P, Krishnatry R, Shago M, et al. BRAF Mutation and CDKN2A Deletion Define a Clinically Distinct Subgroup of Childhood Secondary High-Grade Glioma. *J Clin Oncol* 2015;33:1015–22. doi:10.1200/JCO.2014.58.3922.
- [33] Zhao D, Lu X, Wang G, Lan Z, Liao W, Li J, et al. Synthetic essentiality of chromatin remodelling factor CHD1 in PTEN-deficient cancer. *Nature* 2017;542:484–8. doi:10.1038/nature21357.
- [34] Wang X, Cao X, Sun R, Tang C, Tzankov A, Zhang J, et al. Clinical Significance of PTEN Deletion, Mutation, and Loss of PTEN Expression in De Novo Diffuse Large B-Cell Lymphoma. *Neoplasia* 2018;20:574–93. doi:10.1016/j.neo.2018.03.002.
- [35] Robertson AG, Kim J, Al-Ahmadie H, Bellmunt J, Guo G, Cherniack AD, et al. Comprehensive Molecular Characterization of Muscle-Invasive Bladder Cancer. *Cell* 2017;171:540–556.e25. doi:10.1016/j.cell.2017.09.007.
- [36] Fawcett T. An introduction to ROC analysis. *Pattern Recognit Lett* 2006;27:861–74. doi:10.1016/j.patrec.2005.10.010.
- [37] Friedman JH. Greedy Function Approximation: A Gradient Boosting Machine. *Ann Stat* 2001;29:1189–232.
- [38] Cole AJ, Dwight T, Gill AJ, Dickson K-A, Zhu Y, Clarkson A, et al. Assessing mutant p53 in primary high-grade serous ovarian cancer using immunohistochemistry and massively parallel sequencing. *Sci Rep* 2016;6:26191. doi:10.1038/srep26191.
- [39] Ciconas J, Kvederaviciute K, Meskinyte I, Meskinyte-Kausiliene E, Skeberdyte A, Ciconas J. KRAS, TP53, CDKN2A, SMAD4, BRCA1, and BRCA2 Mutations in Pancreatic Cancer. *Cancers* 2017;9. doi:10.3390/cancers9050042.
- [40] Prior IA, Lewis PD, Mattos C. A comprehensive survey of Ras mutations in cancer. *Cancer Res* 2012;72:2457–67. doi:10.1158/0008-5472.CAN-11-2612.
- [41] Futreal PA, Coin L, Marshall M, Down T, Hubbard T, Wooster R, et al. A census of human cancer genes. *Nat Rev Cancer* 2004;4:177–83. doi:10.1038/nrc1299.
- [42] Sager R. Tumor suppressor genes: the puzzle and the promise. *Science* 1989;246:1406–12. doi:10.1126/science.2574499.
- [43] Irizarry RA, Wu Z, Jaffee HA. Comparison of Affymetrix GeneChip expression measures. *Bioinformatics* 2006;22:789–94. doi:10.1093/bioinformatics/btk046.
- [44] Langfelder P, Horvath S. WGCNA: an R package for weighted correlation network analysis. *BMC Bioinformatics* 2008;9:559. doi:10.1186/1471-2105-9-559.
- [45] Lachmann A, Giorgi FM, Lopez G, Califano A. ARACNe-AP: gene network reverse engineering through adaptive partitioning inference of mutual information. *Bioinforma Oxf Engl* 2016;32:2233–5. doi:10.1093/bioinformatics/btw216.
- [46] Pollen AA, Nowakowski TJ, Shuga J, Wang X, Leyrat AA, Lui JH, et al. Low-coverage single-cell mRNA sequencing reveals cellular heterogeneity and activated signaling pathways in developing cerebral cortex. *Nat Biotechnol* 2014;32:1053–8. doi:10.1038/nbt.2967.
- [47] Risso D, Perraudeau F, Gribkova S, Dudoit S, Vert J-P. A general and flexible method for signal extraction from single-cell RNA-seq data. *Nat Commun* 2018;9. doi:10.1038/s41467-017-02554-5.

- [48] Bush EC, Ray F, Alvarez MJ, Realubit R, Li H, Karan C, et al. PLATE-Seq for genome-wide regulatory network analysis of high-throughput screens | Nature Communications. Nat Commun 2017;8. doi:doi.org/10.1038/s41467-017-00136-z.
- [49] Giorgi FM, Del Fabbro C, Licausi F. Comparative study of RNA-seq- and microarray-derived coexpression networks in *Arabidopsis thaliana*. Bioinforma Oxf Engl 2013;29:717–24. doi:10.1093/bioinformatics/btt053.
- [50] Lawrence M, Huber W, Pagès H, Aboyoun P, Carlson M, Gentleman R, et al. Software for Computing and Annotating Genomic Ranges. PLOS Comput Biol 2013;9:e1003118. doi:10.1371/journal.pcbi.1003118.
- [51] Ashburner M, Ball CA, Blake JA, Botstein D, Butler H, Cherry JM, et al. Gene Ontology: tool for the unification of biology. Nat Genet 2000;25:25–9. doi:10.1038/75556.
- [52] Maaten L van der, Hinton G. Visualizing Data using t-SNE. J Mach Learn Res 2008;9:2579–605.
- [53] Duong T. ks: Kernel Density Estimation and Kernel Discriminant Analysis for Multivariate Data in R. J Stat Softw 2007;021.
- [54] Gautier L, Cope L, Bolstad BM, Irizarry RA. affy—analysis of Affymetrix GeneChip data at the probe level. Bioinformatics 2004;20:307–15. doi:10.1093/bioinformatics/btg405.

## Supporting Information Legends

**Figure S1.** Table of most somatically mutated genes across TCGA tumor samples, in terms of number of samples where the gene is somatically mutated with altered protein product sequence. This table includes also MutSig-blacklisted genes (in grey) such as Titin (TTN), Obscurin (OBSCN) and Mucin genes.

**Figure S2.** Relationship between alteration models and alteration frequency in the Pan-cancer dataset, for mutations (left), amplifications (center) and deletions (right).

**Figure S3.** Performance of Pan-cancer alterations models globally (left) and for MutSig genes, COSMIC oncogenes and COSMIC tumor suppressors. Asterisks indicate a significant (<0.01) difference between a distribution and the global “Other Genes” distribution according to Two-tailed Wilcoxon tests.

**Figure S4.** ROC curves for gbm TP53 models in Breast Cancer, using original expression data, VIPER aggregation (TF “activity”) and WGCNA aggregation (robust tukey biweight average of clusters).

**Figure S5.** AUROCs of EGFR amplication gbm prediction models in Glioblastoma with increasing noise, calculated using gene expression (black line) or aggregated gene expression using the WGCNA (green line) or VIPER (red line) algorithms.

**Figure S6.** AUROCs of EGFR amplication gbm prediction models in Lung Squamous Carcinoma (LUSC) with increasing noise, calculated using gene expression (black line) or aggregated gene expression using the WGCNA (green line) or VIPER (red line) algorithms.

**Figure S7.** AUROCs of PVT1 amplication gbm prediction models in Ovarian Cancer with increasing noise, calculated using gene expression (black line) or aggregated gene expression using the WGCNA (green line) or VIPER (red line) algorithms.

**Figure S8.** AUROCs of PTEN deletion gbm prediction models in Sarcoma with increasing noise, calculated using gene expression (black line) or aggregated gene expression using the WGCNA (green line) or VIPER (red line) algorithms.

**Figure S9.** Beta distribution used to down-sample the 43.8M reads breast cancer sample TCGA-A1-A0SB-01 to 10M reads. The grey line shows the ratio between the target coverage and the original coverage

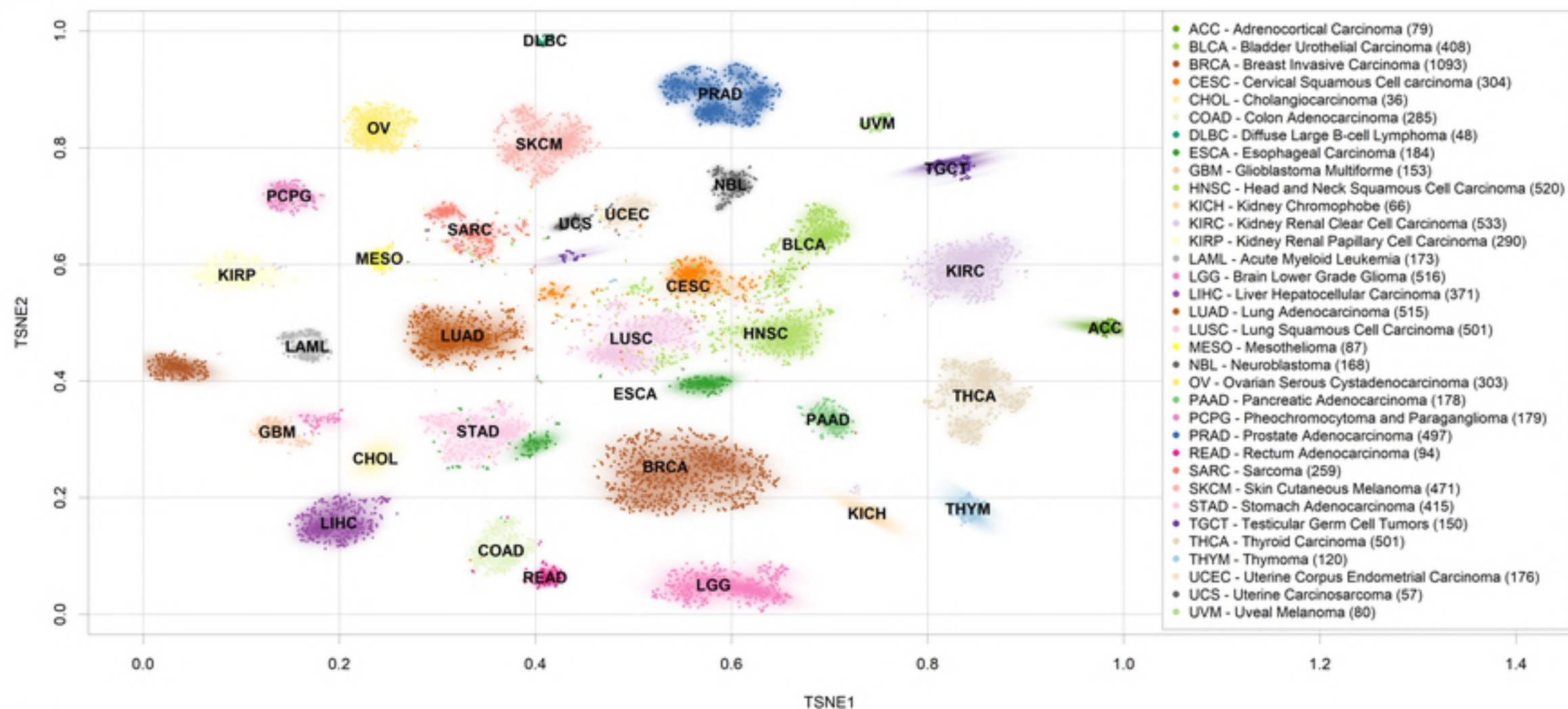
**Supplementary Table S1.** AUROCs for each event in the Pan-Cancer TCGA dataset (24 tumor types with at least 100 samples with co-measured genomic and expression data. The Sheet name indicates the tumor type and genomic alteration type (mut: somatic mutation, amp: amplification, del: deletion).

# 1 **Supplementary Code.** R and bash code snippets used in this study.



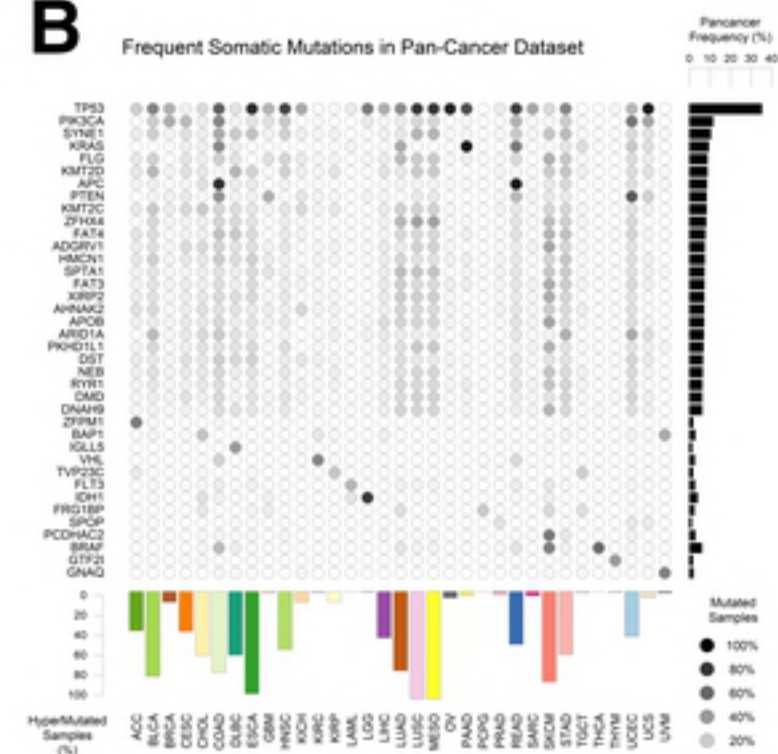
A

## Expression-clustered pancancer dataset



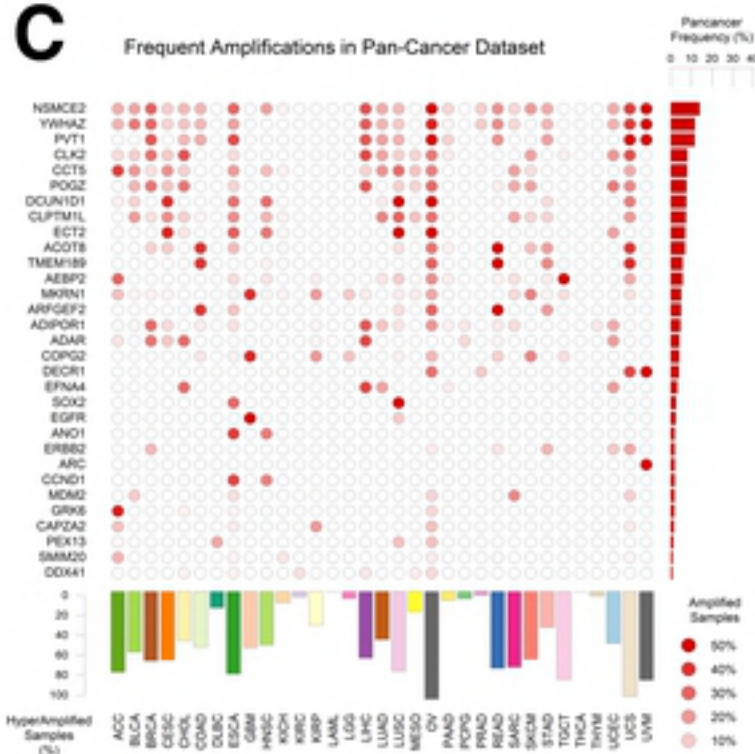
B

## Frequent Somatic Mutations in Pan-Cancer Dataset



C

## Frequent Amplifications in Pan-Cancer Dataset



D

## Frequent Deletions in Pan-Cancer Dataset

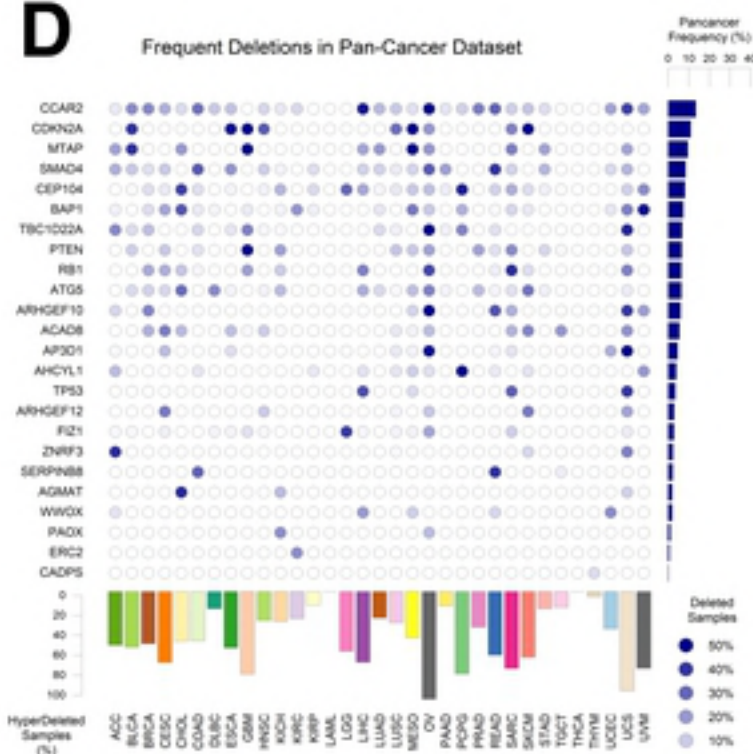


Figure1



## Gene Mutation Classification in Bladder Cancer

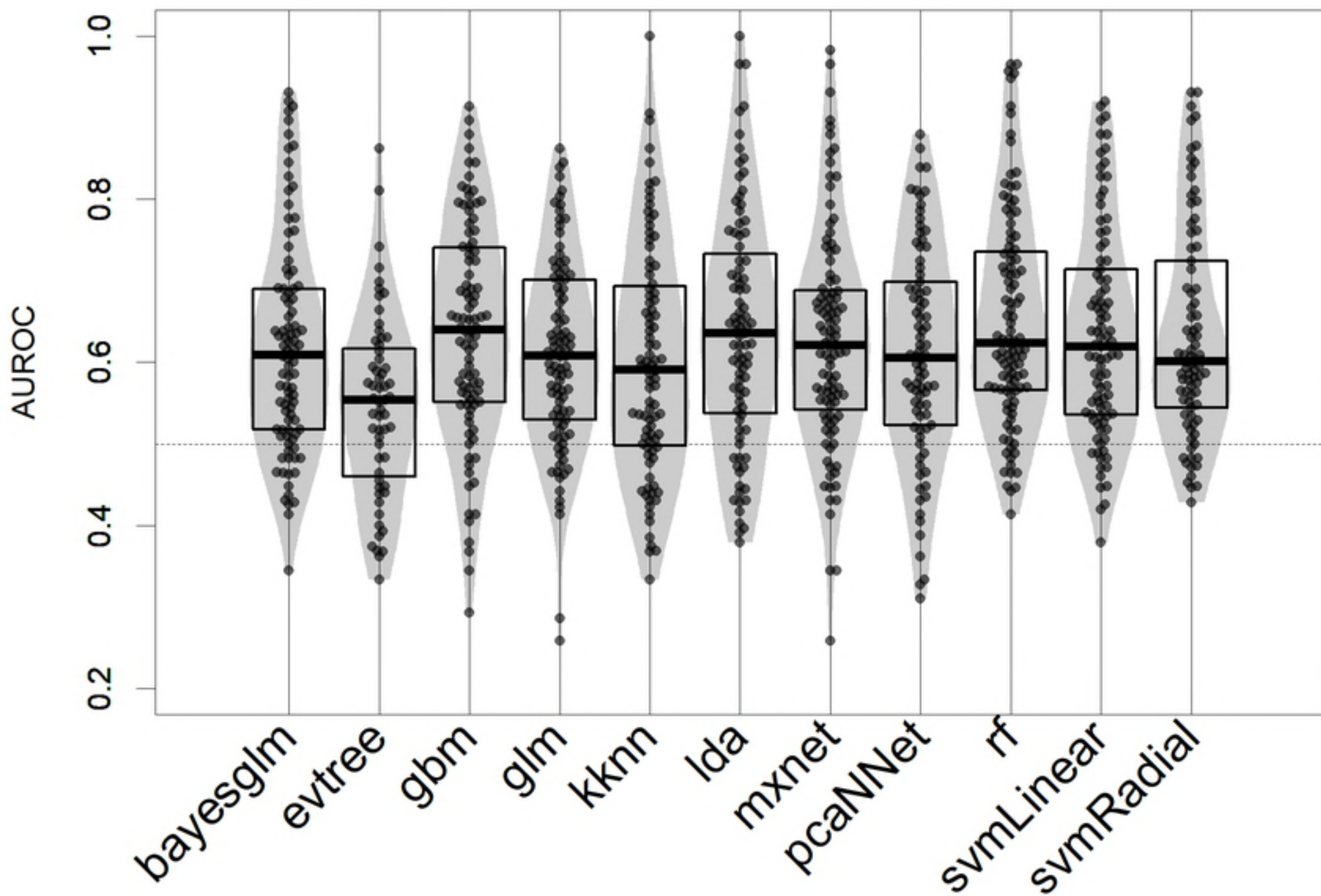


Figure2

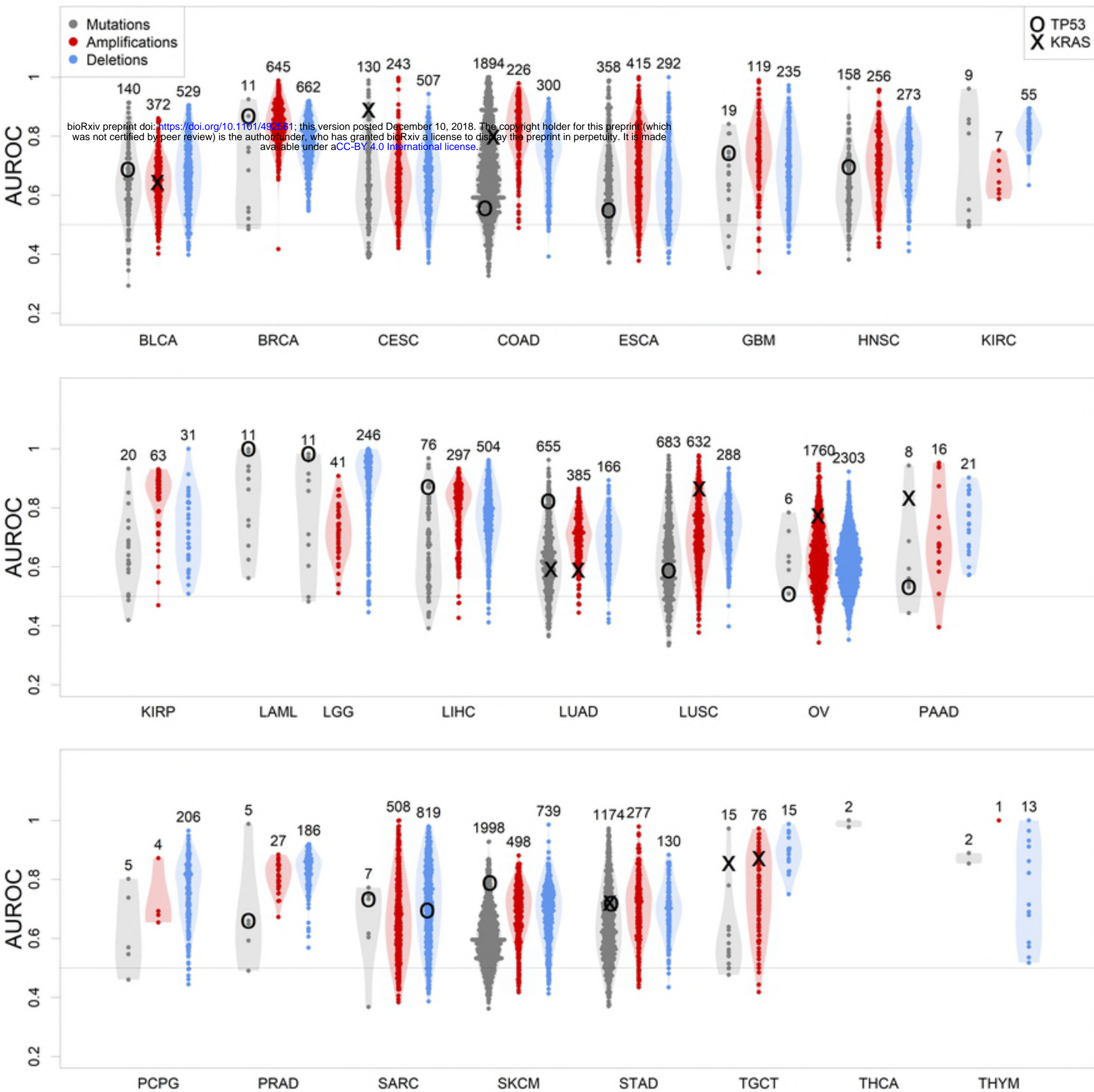
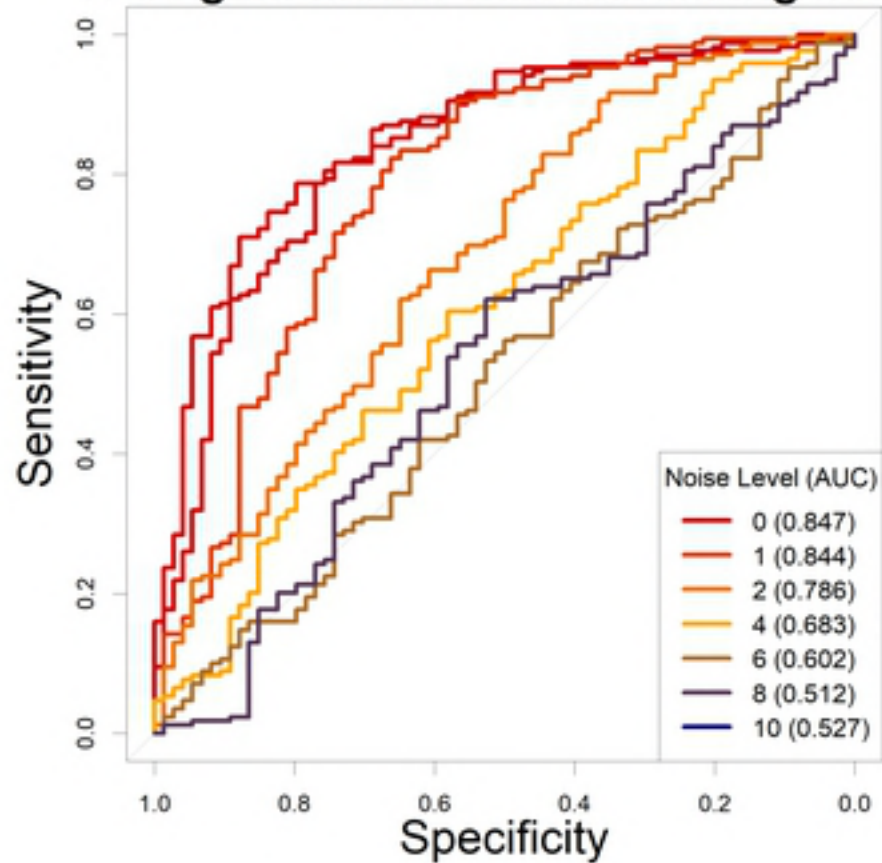
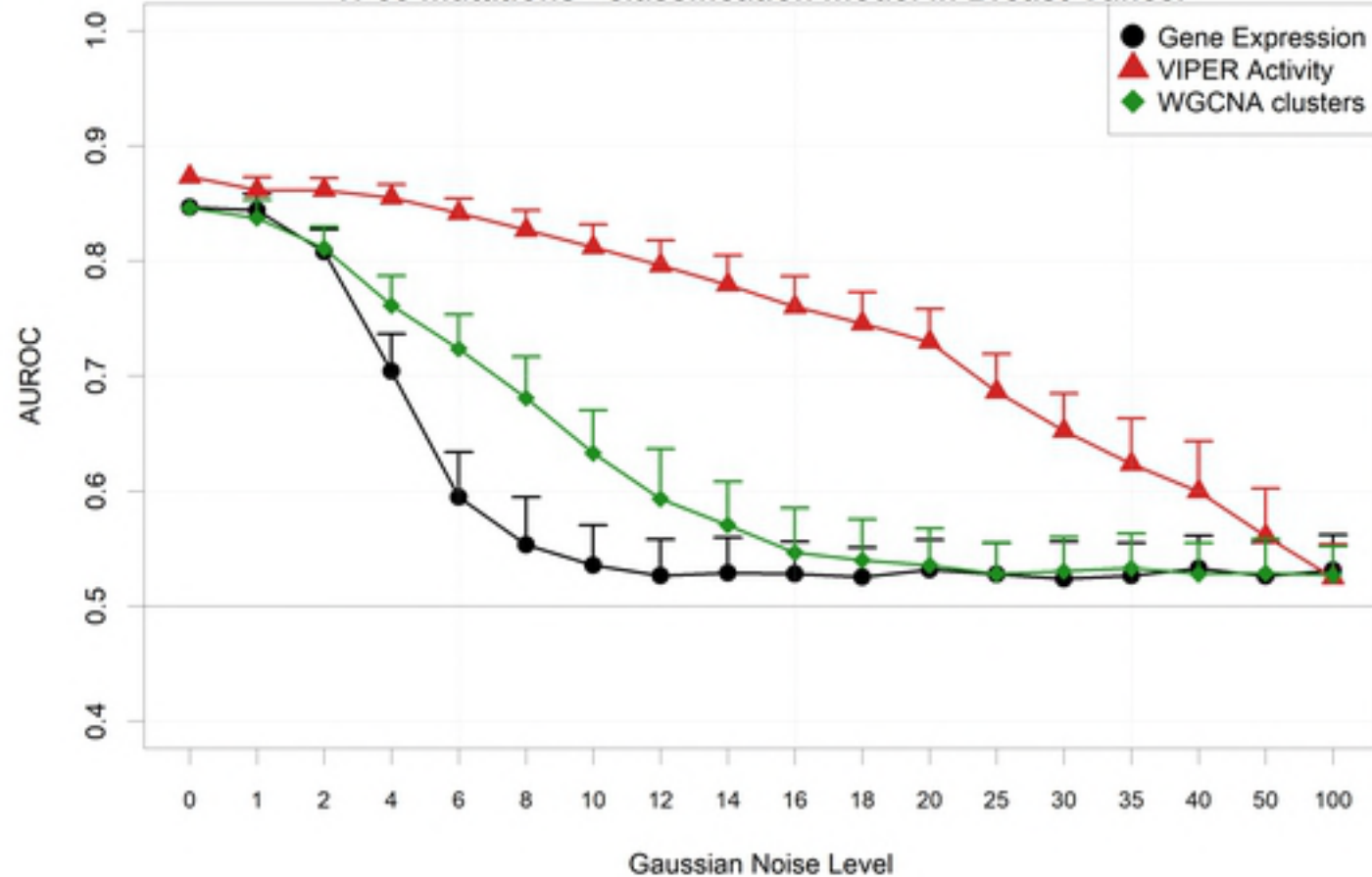


Figure3

**A****TP53 gbmm model with increasing noise****B****TP53 mutations - classification model in Breast Cancer****Figure4**



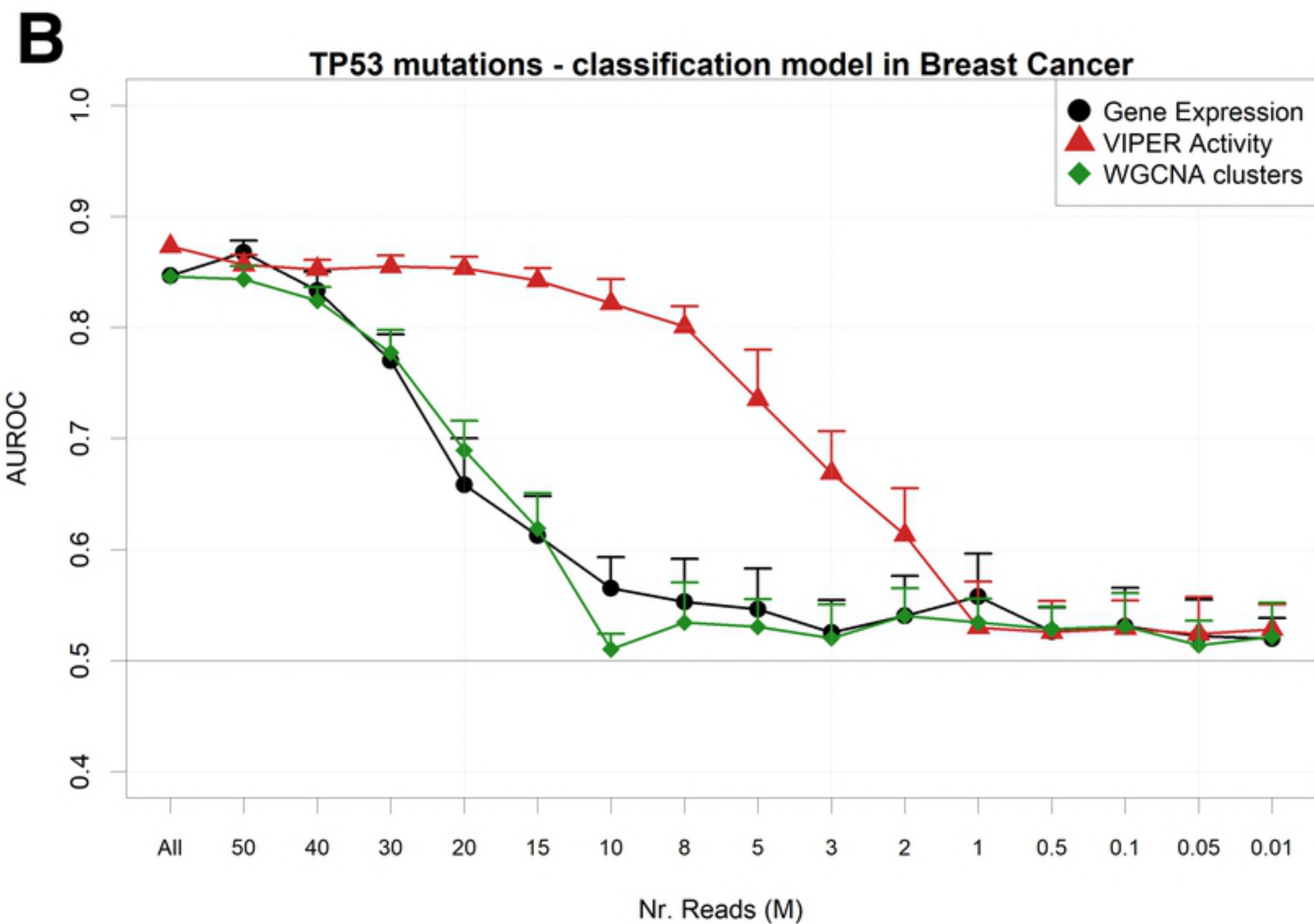
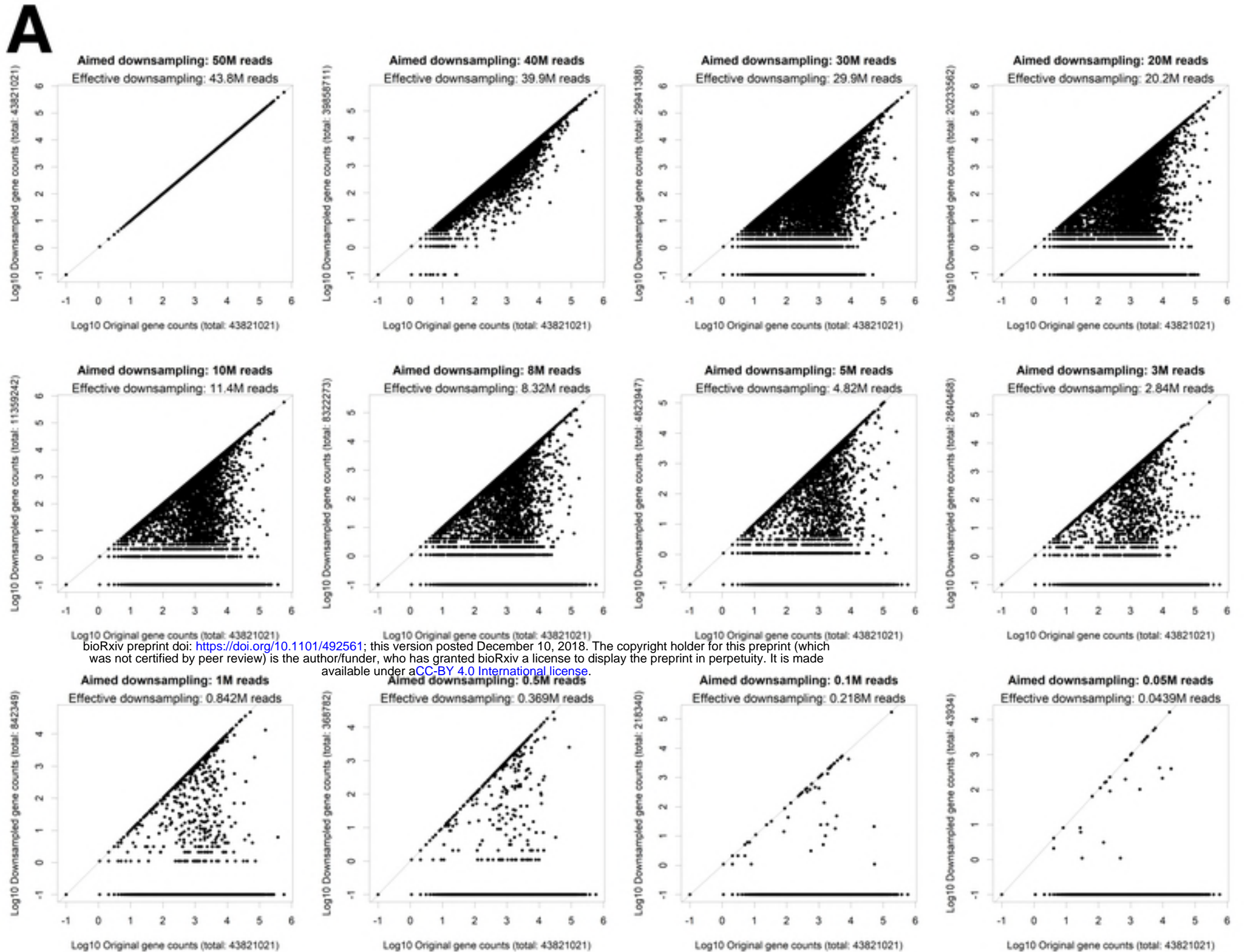


Figure5

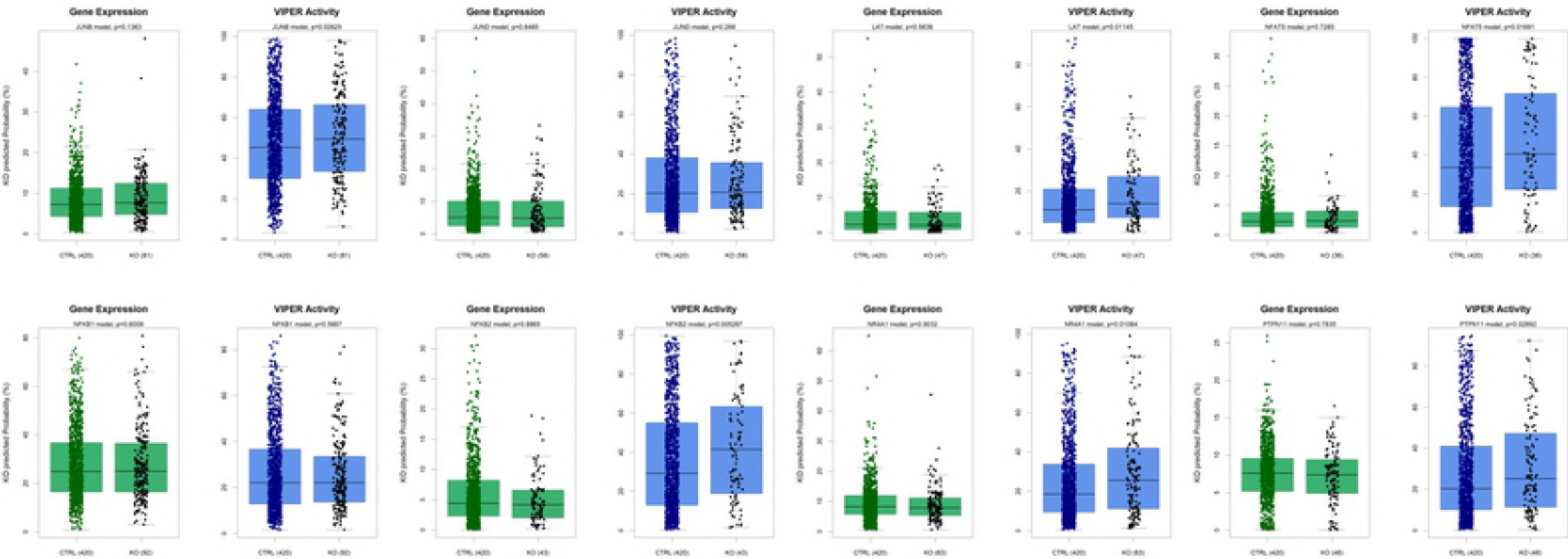


Figure6