

# **MAGICTRICKS: A tool for predicting transcription factors and cofactors that drive gene lists.**

**Roopra, A**

**Dept. of Neuroscience  
5507 WIMR  
1111 Highland Avenue,  
University of Wisconsin-Madison  
Madison,  
WI 53705, USA  
Phone: 1 (608) 265 9072  
asroopra@wisc.edu**

## **ABSTRACT.**

Transcriptomic profiling is an immensely powerful hypothesis generating tool. However, accurately predicting the transcription factors (TFs) and cofactors that drive transcriptomic differences between samples represents a challenge and current approaches are limited by high false discovery rates. This is due to the use of TF binding sequence motifs that, due their small size, are found randomly throughout the genome, and do not allow discovery of cofactors. A second limitation is that even the most advanced approaches that use ChIPseq tracks hosted at sites such as the Encyclopedia Of DNA Elements (ENCODE) assign TFs and cofactors to genes via a binary designation of 'target', or 'non-target' that ignores the intricacies of the biology behind transcriptional regulation.

ENCODE archives ChIPseq tracks of 169 TFs and cofactors assayed in 91 cell lines. The algorithm presented herein, **Mining Gene Cohorts for Transcriptional Regulators Inferred by Kolmogorov-Smirnov Statistics (MAGICTRICKS)**, uses ENCODE ChIPseq data to look for statistical enrichment of TFs and cofactors in gene bodies and flanking regions in gene lists. When compared to 2 commonly used web resources, o-Possum and Enrichr, MAGICTRICKS was able to more accurately predict TFs and cofactors that drive gene changes in 3 settings: 1) A cell line expressing or lacking single TF, 2) Breast tumors divided along PAM50 designations and 3) Whole brain samples from WT mice or mice lacking a single TF in a particular neuronal subtype. In summary, MAGICTRICKS is a standalone application that runs on OSX and Windows machines that produces meaningful predictions of which TFs and cofactors are enriched in a gene list.

## Introduction.

Key to the control of gene expression is the level of transcript in the cell. Transcription factors (TFs) are DNA binding proteins that recognize specific sequence elements (motifs) associated with gene promoters and enhancers. TFs recruit cofactors that do not themselves bind DNA but are brought to promoters via TFs to either enhance or repress gene expression. TFs and cofactors (here on termed 'Factors') are thus key regulators of transcript levels.

It is now routine to obtain the levels of every gene transcript in the genome i.e. whole transcriptome data. The datasets contain tens of thousands of expression values per sample and the number of samples can be in the thousands such as breast cancer transcriptomes archived at The Cancer Genome Atlas (TCGA)<sup>1</sup>.

When comparing transcriptomes from two or more conditions such as normal to cancerous tissue, thousands of mRNA levels can change. We posit that in many cases, the majority of those changes are driven by alterations in the function of a few Factors that coordinate programmatic gene changes on a genome wide scale. Identifying those driving Factors is a fundamental problem and therapeutic opportunity, yet current tools are extremely poor at making such predictions.

Current algorithms fall into two categories. The first looks for motifs that are over-represented in the upstream regions of a list of genes with altered expression compared to the gene population. This approach often fails because with very few exceptions, transcription factor motifs are small (around 6bp) and often have redundant bases. A random 6bp sequence occurs approximately 750,000 times in the human genome. Thus, searches will yield a high false positive rate; most predicted sites will not be functional. Moreover, in many genes, TFs act in concert such that the presence of a single motif element is uninformative. Motif searching also precludes any prediction of cofactors because cofactors do not bind DNA directly. This is especially problematic for translational research where uncovering drugable targets is a major goal, and cofactors are often enzymes that can make excellent drug targets. The second approach relies on assigning genes to Factors (TFs and cofactors) by looking at ChIPSeq tracks such as those archived at ENCODE<sup>1</sup>. Genes are ranked by ChIP signal and then an arbitrary number of 'top' genes are taken as targets. This approach is exemplified by Enrichr<sup>2</sup>. Enrichr is a very powerful tool with excellent performance that allows for assignment of cofactors. However, in general, the binary assignment of genes as targets or non-targets of a Factor does not take into account the biology of Factor regulation of genes. Many Factors have highly non-gaussian ChIP signal distributions; the distributions can have large numbers of sites with very low but detectable signals and skewed, long tails of decreasing ChIP signal. This thwarts unbiased attempts to define genes as 'targets' or 'non-targets'<sup>3</sup>. A new method is required that takes into account the biology of transcriptional regulation in a statistically rigorous manner.

We have devised a novel algorithm that meshes whole transcriptome data with whole genome factor binding (ChIP Seq) data archived at ENCODE. The algorithm is termed ***Mining Gene Cohorts for Transcriptional Regulators Inferred by Kolmogorov-Smirnov Statistics*** (MAGICTRICKS). MAGICTRICKS circumvents the principal confounds of current methods to identify Factors, namely: 1) unacceptably high false positive rates due to the use of TF motif searches 2) inability to identify cofactors due the absence of any binding site motifs 3) arbitrary assignment of Factors to genes based on hard cutoffs of ChIPSeq signals. MAGICTRICKS accepts an input list of genes and compares ChIP signals for the list with the population of ChIP signals for each Factor in ENCODE. We compared MAGICTRICKS to 2 commonly used TF prediction algorithms, oPOSSUM-3<sup>4</sup> and Enrichr<sup>2</sup> using whole transcriptome data from a REST knockdown cell line, TCGA breast tumors and mouse brain tissue. In all cases, MAGICTRICKS predicted biologically meaningful TFs and associated cofactors.

## Methods

### MCF7 RNA profiling.

MCF7 cells harboring a REST shRNA knockdown and controls were generated as described<sup>5</sup>. RNA was extracted in triplicate from shControl and shREST cells using TRIZOL, checked on an Agilent NanoChip and subjected to microarray hybridization onto Nimblegen arrays (2006-08-03\_HG18\_60mer\_expr) at the University of Wisconsin-Madison Biotechnology Center. Spot intensities were RMA normalized<sup>6</sup>. Probes were considered present in the experiment if they were 'present' in all 3 samples of either controls or knockdown cells. If a probe was present in only one condition, the values in the other condition were accepted regardless of the present/absent call. Probes were then collapsed to a single gene value per sample by taking the median value of all probes per gene per sample. Fold changes between control and REST knock-down cells were calculated along with Benjamini-Hochberg corrected Student t-test p values. 'Up-regulated' and 'down-regulated' gene lists for MAGICTRICKS analysis were generated by taking any genes changed more than 3 Standard Deviations from the mean fold change and having a corrected p value < 0.05.

### TCGA profiling.

Breast cancer RNAseq data (data\_RNA\_Seq\_v2\_expression\_median.txt) was downloaded from cBioPortal (latest datafreeze as of 6July2017). PAM50<sup>7</sup> designations were used to parse tumors into Luminal-A or Basal subtypes (supplemental File 9 – list of samples and PAM50 designation). A gene with an RPKM  $\geq 1$  in at least 50% of either basal or Luminal-A samples was considered for further analysis. Fold changes, statistics and gene lists were then generated as above for REST knock-down cells.

### Sams et al profiling.

RNAseq data (GSE84175) was downloaded from GEO. A gene with an RPKM  $\geq 1$  in all controls or all CTCF KO cells was considered for further analysis and fold changes, statistics and gene lists were then generated as above for REST knock-down cells.

### MAGICTRICKS.

MAGICTRICKS aims to determine whether genes in a query list are associated with higher ChIP values than expected by chance for a given transcription factor or cofactor based on ChIPseq tracks archived at ENCODE. To do this, we assigned a ChIP value for each factor in ENCODE at every gene. We then generated an algorithm to test the hypothesis that a list of query genes are enriched for high ChIP signals.

### Generation of the MAGIC Matrix.

We defined a single ChIP signal for each factor in ENCODE at every gene as follows:

All genes in the human genome (hg19) were assigned a gene domain that was defined as 5Kb either side of the gene body (5'UTR – 5Kb to 3'UTR + 5Kb). NarrowPeak files for all transcription factors and cofactors were downloaded from <http://genome.ucsc.edu/cgi-bin/hgFileUi?db=hg19&q=wgEncodeAwgTfbsUniform>. These files provide called peaks of signal enrichment based on pooled, normalized (interpreted) data. Custom Python scripts were then used to extract the highest ChIPseq peak value (signalValue) for a given factor seen in any cell line within every gene domain. Thus each factor and gene were assigned a single signalValue – the highest value – observed in ENCODE. The highest value was chosen to 1) maximize the likelihood that this site will be bound in cell types not present in ENCODE and 2) Work by others suggest that the strongest sites may be the most functional<sup>3</sup>. The resulting array (factor columns vs gene rows) was normalized such that the sum of all signalValues for any factor (columns) was arbitrarily set at 100,000. We term the final array the MAGIC Matrix (Supplemental File 1).

### Defining Enriched Factors in a Gene list.

Two lists are entered into MAGICTRICKS. The first is a master list containing all genes of relevance to the experiment (for example, all genes deemed to be present or expressed in an experiment). The second list contains genes of interest – the query list (e.g. genes induced above a threshold in one condition over control).

Query lists are subsets of the master list and both lists are filtered by MAGICTRICKS for genes present in the MAGIC matrix.

For each Factor in the Matrix, MAGICTRICKS orders the N genes in the accepted master list by ChIP signal (lowest to highest) obtained from MAGIC Matrix. It then generates a normalized cumulative distribution function (CDF) such that:

$$M(c) = \frac{1}{N} \sum_{i=1}^N 1_{M_i \leq c}$$

Where  $M(c)$  is the fraction of genes in the Master list (N genes) with ChIP signal less than c, and M are the ordered ChIP signals.  $1_{M_i \leq c}$  is an indicator function that is equal to 1 if  $M_i \leq c$  else equal to 0. An empirical distribution function,  $Q(c)$  is generated by similarly ordering the X query genes Q, where  $Q \subset M$ , and  $1_{Q_i \leq c} = 1$  if  $Q_i \leq c$  else  $1_{Q_i \leq c} = 0$ :

$$Q(c) = \frac{1}{X} \sum_{i=1}^X 1_{Q_i \leq c}$$

It then screens for Factors where the query distribution is right-shifted compared to the master CDF. This is accomplished by calculating the difference between the 2 distributions and defining the supremum ( $D_S$ ) and infimum  $D_I$  of the difference as:

$$D_S = \sup_c \{M(c) - Q(c)\}$$

$$D_I = \inf_c \{M(c) - Q(c)\}$$

If the cumulative of query samples is left shifted compared to the population ( $|D_S| < |D_I|$ ), the factor is triaged and not considered further. If the query cumulative is right shifted compared to the population cumulative i.e.  $|D_S| > |D_I|$ , MAGICTRICKS performs a 1-tailed Kolmogorov-Smirnov (KS) test between the master and query list to test the null hypothesis that the query list is drawn randomly from the master;  $D_S$  equals the KS statistic. The argument of the KS statistic is the 'critical ChIP' and represents the minimum ChIP value a gene must have to be considered a target gene for the Factor.

A Score for each Factor is calculated to incorporate the Benjamini-Hochberg corrected KS p value, as well as a measure of how the highest ChIP values in the query list compare to the highest values in the population. Thus, for each Factor, the mean of the top n chip signals in the query list is compared to the mean of the top n ChIPs in the population to produce a ratio where:

$$n = 0.05X$$

and

$$r = \frac{\mu_q}{\mu_m}$$

Where  $\mu_q$  is the mean for the top n signals for Q and  $\mu_m$  is the mean of the top n signals for M.

A Score (S) is assigned to each Factor thus:

$$S = -\log(P_{corr}) \times r$$

Where  $P_{corr}$  is the Benjamini-Hochberg corrected KS 1-tailed p value. Factors are then sorted by Score.

## Implementation of MAGICTRICKS.

A tab or comma delimited text file is requested by MAGICTRICKS (lists file). The first column contains a list of all genes of relevance to the experiment (master list). This could be a list of all genes expressed in the system for example, or all genes with a detectable signal. Any number of other columns are then added that contain query lists. For example, a query list may be all genes that go up under some criterion and another may be all genes that go down. The first row is the header and must have unique names for each column.

Using the algorithm above, MAGICTRICKS analyzes each query list and generate a series of output files and sub-directories in the directory containing the lists file. An 'Accepted\_Lists.txt' file is generated which is the original lists file filtered for genes in the MAGIC Matrix. A Platform\_Matrix.txt file is the MAGIC Matrix filtered for genes in the the master list. A series of sub-directories are generated named after each query list in the lists file. Each directory contains sub-directories and files for the stand-alone analysis of that query list. A query\_list\_summary.pdf contains a bar graph of Factors and Scores with  $P_{corr} < 0.05$  (e.g. Fig. 1b). Query\_list\_Summary.xls is a spreadsheet containing statistical information for all non-triaged factors. Data reported in Query\_list\_Summary.xls are:

Critical ChIP:	ChIP value at $d_{sup}$ (This value is used to determine target genes in the list)
Obs Tail Mean:	Average of the 5 <sup>th</sup> percentile ChIP values (n values) in the query list.
Exp Tail Mean:	Average of the top n ChIP values in the population.
Tail Enrichment:	Ratio of the Obs and Exp Tail Means (r)
Raw P:	KS p value
Corrected P (FDR):	Benjamini Hochberg corrected p value ( $P_{corr}$ )
Score:	Score ( $-\log(P_{corr}) \times r$ )

All Factors associated with  $P_{corr} < 0.05$  are highlighted in bold red.

Query\_list\_Targets.gmx is a tab delimited file in the GMX format utilized by GSEA<sup>8</sup>. Columns contain all target genes of a Factor whose ChIP signal is greater than the Critical ChIP i.e. the ChIP at which there is the maximal difference between the population and query cumulative. The GMX format will allow this file to be used directly in any subsequent GSEA analysis. The first column contains all genes in the analysis. The second column contains genes in the input list. All subsequent columns contain target genes for the Factor that is named at the top of the column.

A sub-directory called 'CDFs' contains graphical displays of the analysis for all factors with  $P_{corr} < 0.05$ , named 'rank'\_factor'.pdf (e.g. 1\_NRSF.pdf; rank = 1, factor = NRSF) (Fig. 1c) where ranking is determined by Score. Two cumulative functions are displayed: the black curve is the fractional cumulative of all genes in the master list against the ChIP values, red is the same for query genes. A blue vertical line denotes the Critical ChIP value i.e. ChIP value at  $d_{sup}$ . Red ticks represent each gene in the query list and black ticks are all genes in the population. Red ticks with circles ('lollipops') are those n query genes in the 5<sup>th</sup> percentile of ChIP values. Black lollipops are genes in the population with the n highest ChIP values.

A second sub-directory called 'Auxiliary\_Files' is populated with data behind the summary files. The 'query\_list\_raw\_results.txt' file contains the same columns as the Query\_list\_Summary.xls file but has raw data for all factors including those that were triaged and not considered for further statistical analysis. It also contains the KS statistic for each factor. KS statistics with a negative sign denote KS values for triaged factors; the negative sign is used by the algorithm for triage sorting.

Query\_list\_Sub\_Matrix.txt is the MAGIC Matrix filtered for genes in the query list.

Triaged\_Factors.txt is a list of factors that were not considered (triated) because  $|d_{inf}| > |d_{sup}|$



Triaged\_Genes.txt contains all genes in the query that were not in MAGIC Matrix and therefore eliminated from analysis.

A sub-directory named 'Target\_Data' contains comma separated text files for each Factor with  $P_{\text{corr}} < 0.05$ . The file contains the list of target genes for that Factor and associated MAGIC Matrix ChIP value. Within 'Target\_Data' a directory named '\_FDR\_above\_5\_percent' contains lists of target genes and associated ChIP values for Factors with an  $P_{\text{corr}} \geq 0.05$ .

## Obtaining MAGICTRICKS

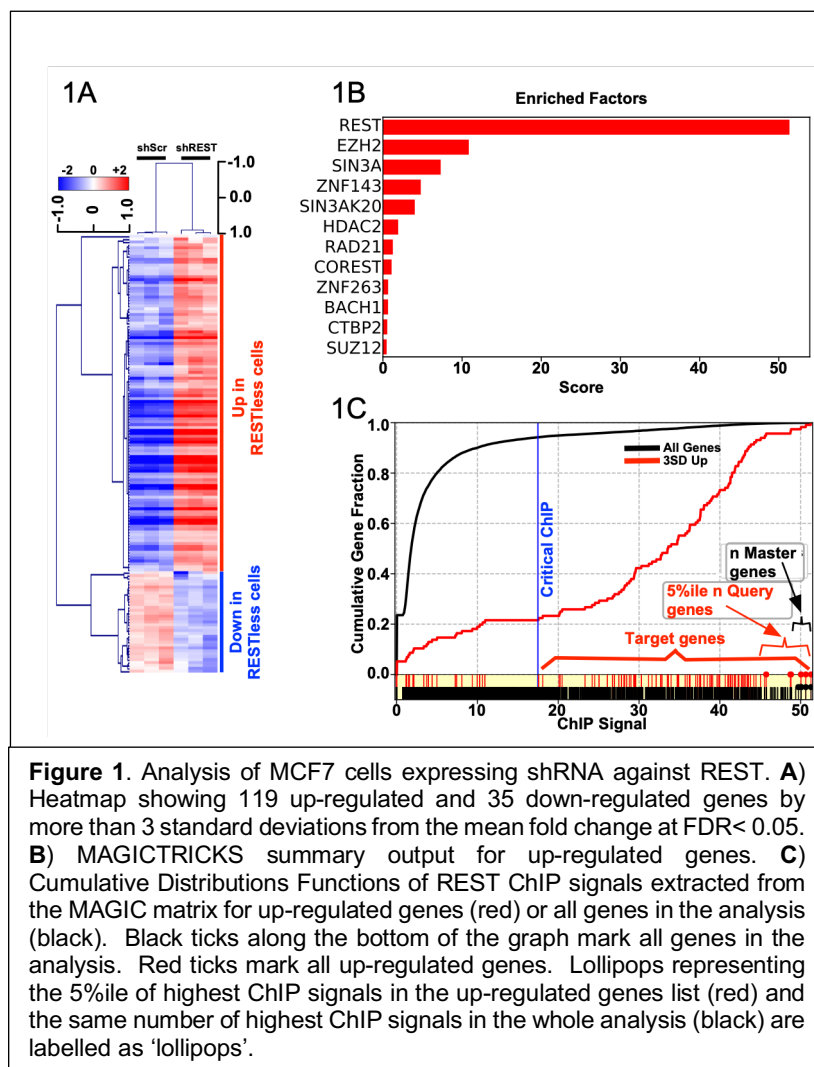
MAGICTRICKS can be downloaded from <https://go.wisc.edu/MAGICTRICKS>.

## Results

The purpose of MAGICTRICKS is to identify transcription factors and cofactors that are responsible for the major patterns of gene expression changes in a disease situation or other biologic perturbation. To test MAGICTRICKS, we asked whether MAGICTRICKS could predict which transcription factor was knocked down in a cell line. We have shown that REST is a tumor suppressor whose loss drives tumor growth of MCF7 cells in mouse xenograft model of breast cancer. We have previously published on the transcriptome of MCF7 cells infected with lentivirus expressing either a scrambled shRNA or shRNA targeting REST<sup>5,9,10</sup>. This simple system consists of a clonal cell line lacking a single factor. RNA was subjected to hybridization to Nimblegen arrays, and all genes were assigned a Fold Change ( $FC = \log_2(\text{shREST}/\text{shScramble})$ ) and those with  $FC > 3$  Standard Deviations from the mean FC and  $FDR < 0.05$  were designated as the 'Upregulated' gene list upon loss of REST (118 genes) or 'Downregulated' (35 genes not including REST itself). These 153 genes were able to distinguish shScramble from shREST

samples in unsupervised hierarchical clustering (Fig.1A). Using 15,445 expressed and 118 upregulated genes (supplemental file 2), MAGICTRICKS predicted REST as the highest scoring Factor associated with upregulated genes in shREST cells (raw  $p = 7.6 \times 10^{-55}$ ,  $FDR = 3 \times 10^{-53}$ , Score = 51.3). MAGICTRICKS also classified the REST corepressors SIN3A, CoREST, HDAC2, and CtBP2 as significant at  $FDR < 0.05$ . (Fig.1B and supplemental file 3) (see discussion). Figure 1C shows the highly divergent cumulatives of the population of REST ChIP values (black) and those associated with genes induced more than 3 standard deviations of the mean change (red).

For comparison, we used the 3SD gene list in oPossum. Although oPossum did predict REST as the most likely driver (z score = 83) (supplemental file 4), oPossum predicted 45 other transcription factors at  $z > 1.65$ . Most of the factors have no known physical or functional association with REST. REST is a unique TF in that it has an extended binding motif of 23bp, which likely aided its identification by oPossum. Regardless, oPossum was not able to call any REST cofactors as drivers of the gene list.



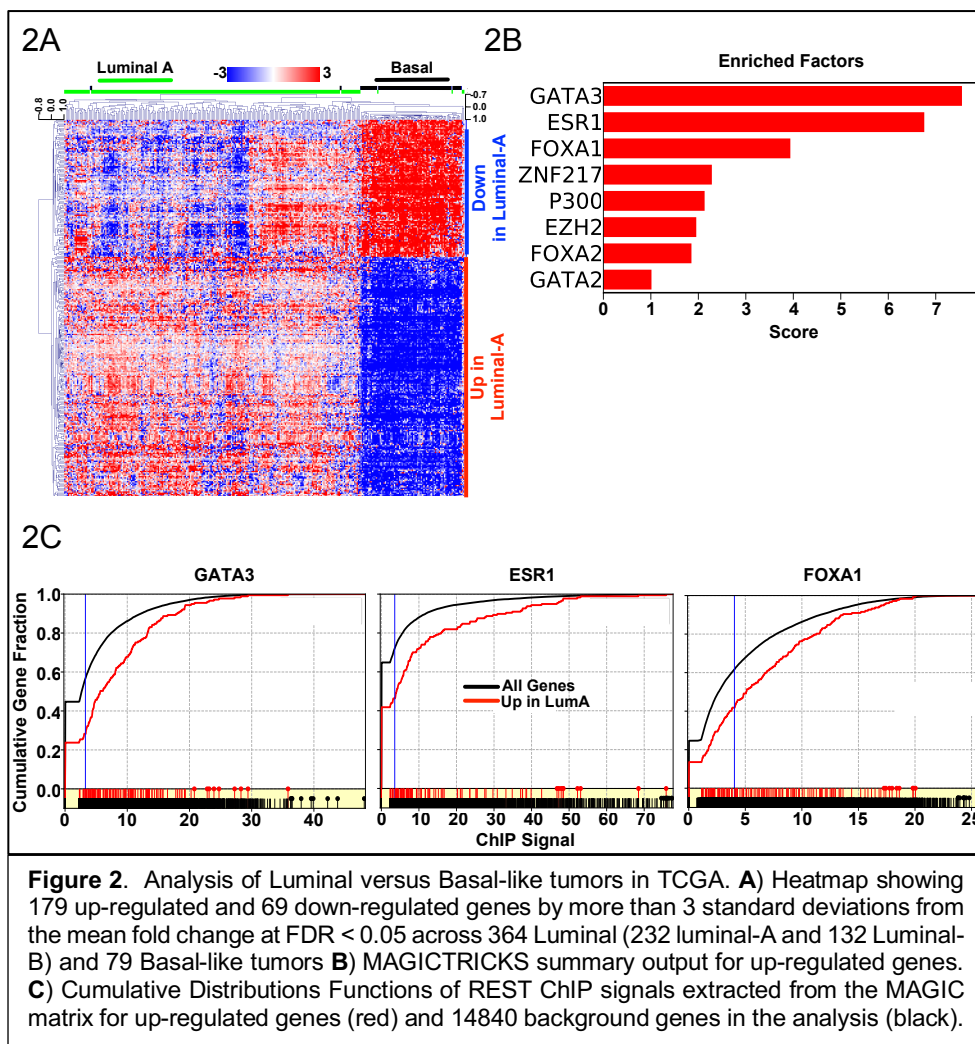
**Figure 1.** Analysis of MCF7 cells expressing shRNA against REST. **A)** Heatmap showing 119 up-regulated and 35 down-regulated genes by more than 3 standard deviations from the mean fold change at  $FDR < 0.05$ . **B)** MAGICTRICKS summary output for up-regulated genes. **C)** Cumulative Distributions Functions of REST ChIP signals extracted from the MAGIC matrix for up-regulated genes (red) or all genes in the analysis (black). Black ticks along the bottom of the graph mark all genes in the analysis. Red ticks mark all up-regulated genes. Lollipops representing the 5%ile of highest ChIP signals in the up-regulated genes list (red) and the same number of highest ChIP signals in the whole analysis (black) are labelled as 'lollipops'.

The same list yielded REST as the top hit in 13 ENCODE cell lines using Enrichr followed by EZH2 and HDAC2. No other Factors were significant at FDR<0.05 (data not shown)

Next we tested MAGICTRICKS on The Cancer Genome Atlas (TCGA) Nature 2012<sup>11</sup> provisional breast cancer dataset. We extracted 364 Luminal and 79 Basal-like tumors. Luminal tumors are a subset of breast cancers defined by their robust expression of estrogen receptor alpha (ER $\alpha$ ) and associated pioneer factors. They also express the progesterone and her2 receptors. Basal-like tumors are estrogen receptor negative and also lack progesterone and her2 receptors. This is a complex dataset with over 400 samples and heterogeneous tissue. We took all genes up-regulated more than 3SD from the mean fold change in Luminal tumors relative to Basal-like (203 genes) as the gene list for input into MAGICTRICKS and 17814 expressed genes (supplemental file 5). MAGICTRICKS called GATA3, ER $\alpha$  (ESR1) and FOXA1 as the highest scoring factors (GATA3:  $p=8.4 \times 10^{-14}$ , FDR =  $2.9 \times 10^{-12}$ , Score = 7.6; ESR1:  $p=1.9 \times 10^{-11}$ , FDR =  $3.4 \times 10^{-10}$ , Score = 6.8; FOXA1:  $p=5.7 \times 10^{-7}$ , FDR =  $6.7 \times 10^{-6}$ , Score = 3.9). These 3

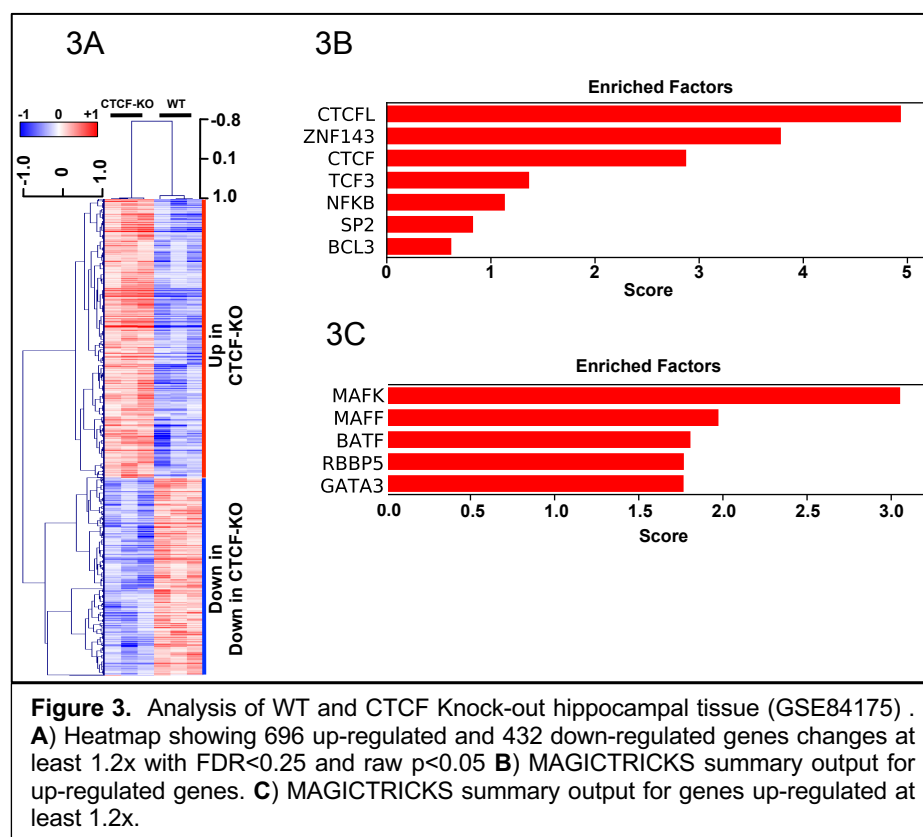
Factors characterize the luminal breast cancer phenotype and can regulate their mutual expression<sup>12-19</sup>. Importantly, MAGICTRICKS also called ER $\alpha$  cofactors EZH2<sup>20</sup>, P300<sup>21</sup> and ZNF217<sup>22</sup> as significant Factors at FDR<0.05 (Fig.2 and supplemental file 6). The same gene list yielded 46 TFs with oPossum at  $z < 1.65$  with ER $\alpha$  ranked at 41 ( $z=2.633$ ). FOXA2 ( $z=23.2$ ) and FOXA1 ( $z=20.1$ ) were the top hits (supplemental file 7). Thus, whereas MAGICTRICKS predicted ER $\alpha$  as a principal factor and identified its cohort of pioneer factors and cofactors in ER positive breast cancer, oPossum was unable to clearly underscore the estrogen receptor biology in this test.

At the same FDR cutoff (<0.05), Enrichr yielded ESR1, GATA3, P300, TCF12, EZH2 and FOXM1 (supplemental file 7 – Enrichr tab). Thus there was greater congruence between MAGICTRICKS and Enrichr. Interestingly, whereas Enrichr provided a list of 20 ESR1 target genes, MAGICTRICKS offered 95 Estrogen Receptor targets as defined by having a ChIP signal greater than the Critical Chip (supplemental file 6 – ESR1 targets tab). To test whether MAGICTRICKS was providing a biologically meaningful expanded target gene list rather than merely a large list of genes unrelated to estrogen receptor biology, we performed ontological analysis on target



genes from Enrichr and MAGICTRICKS. The 20 ESR1 Enrichr and 95 ESR1 MAGICTRICKS targets were entered into STRING<sup>23-25</sup> with the 17814 expressed genes as background. Both lists highlighted 'Estrogen-dependent gene expression' as the top Reactome<sup>26</sup> term. Three of the 20 Enrichr genes were associated with 'Estrogen-dependent gene expression' at FDR=0.0112 whereas 7 of the 95 MAGICTRICKS genes associated at FDR=5x10<sup>-4</sup> (supplemental file 8). Importantly, whereas MAGICTRICKS tagged ESR1, FOXA1 and PGR1 as ESR1 targets consistent with known estrogen receptor regulation<sup>27-29</sup>, Enrichr failed to call these crucial genes as targets. Further, under the STRING 'Biological Processes' output, the MAGICTRICK target list was tagged with the terms 'gland morphogenesis' (FDR=8.5x10<sup>-5</sup>), 'branching involved in mammary gland duct morphogenesis' and 'mammary gland alveolus development' (FDR=4.5x10<sup>-3</sup>) amongst other terms associated with mammary gland development and pregnancy associated remodelling (supplemental file 8). The Enrichr target list was not associated with these terms.

MAGICTRICKS was developed using ChIPseq data derived from immortalized or transformed human cancer cell lines and the above two examples utilize either an immortalized cell line or cancer tissue. To test whether MAGICTRICKS could be used on data derived from non-cancerous, highly heterogeneous, rodent samples, we turned to non-malignant mouse brain tissue. CTCF is a TF that organizes long distance interactions in the genome. Sams et al demonstrated that elimination of CTCF in excitatory neurons (while sparing its expression in other neuron subtypes, astrocytes and glia) results in defects in learning, memory and neuronal plasticity<sup>30</sup>. RNA from whole hippocampus (i.e. all cell types in the formation) from WT and CTCF knock-out mice (n=3 per condition) was subjected to RNAseq and the transcriptome



**Figure 3.** Analysis of WT and CTCF Knock-out hippocampal tissue (GSE84175). **A)** Heatmap showing 696 up-regulated and 432 down-regulated genes changes at least 1.2x with FDR<0.25 and raw p<0.05 **B)** MAGICTRICKS summary output for up-regulated genes. **C)** MAGICTRICKS summary output for genes up-regulated at least 1.2x.

data was archived as GSE84175. We generated two gene lists from this dataset: genes up or down-regulated upon CTCF loss at least 1.2x with FDR<0.25 and raw p<0.05 (696 and 432 genes respectively – supplemental file 10) (Fig. 3A). Figure 3B and supplemental file 11 ('Down' tab) shows that the 3 highest scoring Factors targeting genes that are down-regulated in CTCF-KO hippocampii are CTCFL, ZNF143 and CTCF itself. Interestingly, for the list of up-regulated genes, MAGICTRICKS yielded MAFF and MAFK as top hits (Fig. 3C and supplemental file 11 'Up' tab) which are associated with Locus Control Regions that require CTCF function<sup>31</sup>.

oPOSSUM yielded ZNF354C as the top hit with CTCF ranked number 36 (z=3.8) for down-regulated genes (supplemental file 12 – oPossum tab). Enrichr yielded CTCF as the top hit and was thus more congruent with MAGICTRICKS (supplemental file 1 – Enrichr tab).



In summary, MAGICTRICKS was able to predict which TFs and cofactors drive gene changes in three different types of transcriptomic analysis: a cell line with or without REST, a tumor cohort with or without ER $\alpha$  and brain samples with or without CTCF.

## Discussion.

We introduce MAGICTRICKS as a standalone application that predicts transcription factors and cofactors that drive gene expression changes in transcriptomic experiments. In comparison to a commonly used resource, oPossum, MAGICTRICKS predicted factors that were directly relevant to the biology behind the experiment. MAGICTRICKS predicted both the transcription factors whose binding are enriched in the input list as well as cofactors that are consistent with the known biology of the cognate transcription factor.

MAGICTRICKS was developed in response to an unmet need that would allow allocation of transcription factors and cofactors to lists of genes obtained through whole transcriptome experiments. A widely used approach - exemplified by oPossum - rely on searching for short transcription factor binding motifs in promoter proximal sequences. Using hypergeometric tests provides a likelihood of enrichment for the factor. This approach is limited due to the small size of binding sites for most factors; a random 6bp sequence will be found approximately 750,000 times in the human genome. This approach also cannot predict cofactors as cofactors do not bind DNA directly.

An alternative approach relies on comparing a query list of genes with gene lists attributed to factors based on an arbitrary minimal signal or an arbitrary rank cutoff based on ChIPseq experiments. This approach is exemplified by Enrichr<sup>2</sup> that accepts input lists and predicts factors based on ChIP and ENCODE as well as literature mining. Enrichr works by taking the top 2000 ChIPed genes for each Factor in each cell line in ENCODE as the 'target list' for that Factor. Fisher tests are then performed to assess association between the query list and the 2000 targets. Enrichr was accurately able to predict REST, ESR1 and CTCF in the 3 examples used in this manuscript. The noticeable difference between MAGICTRICKS and Enrichr was in the identification of Factor target genes. Enrichr identifies target genes in the user list by Venn analysis of the user list and pre-defined targets of a Factor; MAGICTRICKS identifies targets by searching for ChIP signals greater than the argument of the Kolmogorov-Smirnov statistic. In the case of ER $\alpha$  when comparing the transcriptomes of Luminal-A and basal-like tumors, ER $\alpha$  targets defined by Enrichr did not include ER $\alpha$  itself, nor PGR or FOXA1 whereas these three key ER $\alpha$  targets<sup>27-29</sup> were in the set of ER $\alpha$  MAGICTRICKS targets. The larger list of ER $\alpha$  targets defined by MAGICTRICKS were enriched for genes associated with estrogen receptor biology as assessed by ontology (Supplemental File 8).

Though Enrichr is a very powerful, userfriendly and informative resource, the arbitrary boundary of target/non-target genes set at 2000 genes does not take into account the fact that many factors have highly skewed binding profiles: many genes will show some, low level of binding/chip signal for a given Factor. For example, extracting all REST ChIPseq tracks from wgEncodeRegTfbsClusteredV2.bed from the UCSC genome browser (at <http://hgdownload.soe.ucsc.edu/goldenPath/hg19/encodeDCC/wgEncodeRegTfbsClustered/>) showed 27,386 binding sites across the human genome with 17,971 genes showing a detectable signal within 5Kb of the promoter (data not shown). This thwarts unbiased attempts to define genes as 'targets' or 'non-targets'. MAGICTRICKS overcomes these shortfalls by utilizing complete ChIPseq profiles without parsing Factors into 'bound' and 'not bound' classes.

Testing MAGICTRICKS on cells transfected with REST shRNA highlights the ability of the algorithm to not only identify transcription factors but also their associated cofactors. MAGICTRICKS called REST as the highest scoring factor but also called SIN3A, HDAC2, CoREST and CtBP2, which are well-characterized REST cofactors<sup>32-35</sup>. EZH2 and SUZ12 were also called as factors that preferentially bind up-regulated genes. Though there are reports of REST interacting with Polycomb<sup>36</sup>, it is likely that Polycomb targets many of the same genes as REST but does so independently<sup>37</sup>. In either case, using a well defined clonal cell line, MAGICTRICKS calls Factors that are either cofactors for REST or co-regulate the same genes. Using gene lists derived from over

300 breast tumor transcriptomes, MAGICTRICKS also correctly called ER $\alpha$  as the principal driver of gene induction in ER positive (Luminal-A) versus negative (Basal-like) tumors. ER $\alpha$  binds chromatin with pioneer factors and MAGICTRICKS successfully called GATA2, GATA3, FOXA1 and FOXA2<sup>19</sup>. P300 and EZH2 are ER $\alpha$  cofactors and were called by MAGICTRICKS<sup>20,21</sup>.

In comparison, oPossum was able to identify REST as the principal factor in the first analysis by using Fisher tests of the input list against lists of genes defined as transcription factor targets by proximity to sequence motifs. As expected, oPossum was unable to call any cofactors due to an absence of binding motifs for this class of nuclear protein. With the breast cancer samples, oPossum ranked ER $\alpha$  41 out of 116 candidate factors. Thus, using oPossum, a naïve researcher studying transcriptomes from Luminal-A and Basal-like tumors would have had difficulty identifying estrogen receptor as a principal factor in the biology of this disease. However, the ER $\alpha$  pioneer factors FOXA1 and FOXA2 were ranked #1 and 2 by oPossum, which would hint at a role for steroid receptors in the disease. Nevertheless, MAGICTRICKS highlighted estrogen receptor biology as a primary driver of transcriptional differences between the 2 groups of tumors.

ENCODE hosts ChIPseq tracks derived from transformed or immortal cell lines and human embryonic stem cells. The MAGIC Matrix incorporates the highest ChIP signal observed across all cell lines within a gene domain for a given factor. This approach was chosen (rather than taking the mean signal for example) because MAGICTRICKS is a discovery platform and taking the maximum value is an attempt to maximize the chance that a particular site may be bound in other systems. A particular concern was that generating a matrix from homogenous cell lines may preclude MAGICTRICKS's use in analysis of transcriptome data derived from in vivo samples that are not immortalized and likely heterogeneous. However, testing of MAGICTRICKS on GSE84175 demonstrates that RNAseq data derived from control and experimental whole brain extract is handled appropriately. In GSE84175, CTCF was deleted in a subset of post-mitotic cells (eliminated from excitatory neurons but retained in all other neuronal subtypes)<sup>30,38</sup>. RNA was extracted from whole hippocampus, which included glia and astrocytes. MAGICTRICKS correctly called CTCF and CTCFL as the factors most likely to drive gene changes in the experiment. Interestingly, ZNF143 was also in the top 3 hits. This is in keeping with the recent findings of Mourad and Cuvier<sup>38</sup> showing that CTCF and ZNF143 coordinate chromatin border domain formation. Thus, MAGICTRICKS was able to call CTCF as well as a known associated factor when provided transcriptome data from a highly heterogeneous tissue where only a small subset of cells were altered in the experiment.

A clear limitation of the current input matrix is that only factors with peaks within 5Kb of the gene body were assigned to a gene: it is clear that a more complete method will take into account distal enhancers that have major roles in global gene regulation. However, despite the absence of distal enhancer information in the MAGIC Matrix, the above three examples show the utility of MAGICTRICKS and MAGIC Matrix using only gene-proximal factor binding. Another limitation is the reliance on a matrix comprised of 161 ENCODE ChIPseq tracks despite there being over 1000 proteins in the human genome that reside in the nucleus. However, the algorithm can accept user defined ChIP seq data and so can evolve as more tracks become available. Users can add their own ChIPseq data to the matrix. To do so, the user would find all ChIPseq peaks within 5Kb of all gene bodies in the MAGIC Matrix and pick the highest to assign to that gene. The sum of all peaks so defined would be normalized to sum to 100,000 and added to the MAGIC matrix as a column.

In summary, MAGICTRICKS is an application for identifying transcription factors and cofactors that preferentially bind lists of genes and provides comprehensive lists of target genes per factor.

## Figure Legends.

**Figure 1.** Analysis of MCF7 cells expressing shRNA against REST. A) Heatmap showing 119 up-regulated and 35 down-regulated genes by more than 3 standard deviations from the mean fold change at FDR < 0.05. B) MAGICTRICKS summary output for up-regulated genes. C) Cumulative Distributions Functions of REST ChIP signals extracted from the MAGIC matrix for up-regulated genes (red) or all genes in the analysis (black). Black ticks along the bottom of the graph mark all genes in the analysis. Red ticks mark all up-regulated genes. Lollipops representing the 95%ile of highest ChIP signals in the up-regulated genes list (red) and the same number of highest ChIP signals in the whole analysis (black) are labelled as 'lollipops'. D) MAGICTRICKS summary file using genes up-regulated 1.2x at FDR<0.05 in cells expressing shRNA against REST.

**Figure 2.** Analysis of Luminal-A and Basal like tumors in TCGA. A) Heatmap showing 203 up-regulated and 143 down-regulated genes by more than 3 standard deviations from the mean fold change at FDR < 0.05 across 233 Luminal-A and 79 Basal-like tumors B) MAGICTRICKS summary output for up-regulated genes. C) Cumulative Distributions Functions of REST ChIP signals extracted from the MAGIC matrix for up-regulated genes (red) and 17814 background genes in the analysis (black).

**Figure 3.** Analysis of WT and CTCF Knock-out hippocampal tissue (GSE84175) . A) Heatmap showing 696 up-regulated and 432 down-regulated genes changes at least 1.2x with FDR<0.25 and raw p<0.05 B) MAGICTRICKS summary output for up-regulated genes. C) MAGICTRICKS summary output for down-regulated genes.

## Conflicts of Interest.

The author has no conflicts of interest to declare.

## Acknowledgements.

I would like to thank Ray Dingledine, Emory University for critical insights into the statistics behind MAGICTRICKS. I thank to Caroline Alexander, John Svaren and Steve Goldstein (University of Wisconsin at Madison) for helpful edits of the manuscript. I would like to thank Queen because 'It's a Kind of Magic'.

# Bibliography.

- 1 Cancer Genome Atlas Research, N., Weinstein, J. N., Collisson, E. A., Mills, G. B., Shaw, K. R., Ozenberger, B. A., Ellrott, K., Shmulevich, I., Sander, C. & Stuart, J. M. The Cancer Genome Atlas Pan-Cancer analysis project. *Nat Genet* **45**, 1113-1120, doi:10.1038/ng.2764 (2013).
- 2 Chen, E. Y., Tan, C. M., Kou, Y., Duan, Q., Wang, Z., Meirelles, G. V., Clark, N. R. & Ma'ayan, A. Enrichr: interactive and collaborative HTML5 gene list enrichment analysis tool. *BMC Bioinformatics* **14**, 128, doi:10.1186/1471-2105-14-128 (2013).
- 3 Biggin, M. D. Animal transcription networks as highly connected, quantitative continua. *Dev Cell* **21**, 611-626, doi:10.1016/j.devcel.2011.09.008 (2011).
- 4 Kwon, A. T., Arenillas, D. J., Worsley Hunt, R. & Wasserman, W. W. oPOSSUM-3: advanced analysis of regulatory motif over-representation across genes or ChIP-Seq datasets. *G3 (Bethesda)* **2**, 987-1002, doi:10.1534/g3.112.003202 (2012).
- 5 Meyer, K., Albaugh, B., Schoenike, B. & Roopra, A. Type 1 Insulin-Like Growth Factor Receptor/Insulin Receptor Substrate 1 Signaling Confers Pathogenic Activity on Breast Tumor Cells Lacking REST. *Mol Cell Biol* **35**, 2991-3004, doi:10.1128/MCB.01149-14 (2015).
- 6 Bolstad, B. M., Irizarry, R. A., Astrand, M. & Speed, T. P. A comparison of normalization methods for high density oligonucleotide array data based on variance and bias. *Bioinformatics* **19**, 185-193 (2003).
- 7 Nielsen, T. O., Parker, J. S., Leung, S., Voduc, D., Ebbert, M., Vickery, T., Davies, S. R., Snider, J., Stijleman, I. J., Reed, J., Cheang, M. C., Mardis, E. R., Perou, C. M., Bernard, P. S. & Ellis, M. J. A comparison of PAM50 intrinsic subtyping with immunohistochemistry and clinical prognostic factors in tamoxifen-treated estrogen receptor-positive breast cancer. *Clin Cancer Res* **16**, 5222-5232, doi:10.1158/1078-0432.CCR-10-1282 (2010).
- 8 Subramanian, A., Tamayo, P., Mootha, V. K., Mukherjee, S., Ebert, B. L., Gillette, M. A., Paulovich, A., Pomeroy, S. L., Golub, T. R., Lander, E. S. & Mesirov, J. P. Gene set enrichment analysis: a knowledge-based approach for interpreting genome-wide expression profiles. *Proc Natl Acad Sci U S A* **102**, 15545-15550, doi:10.1073/pnas.0506580102 (2005).
- 9 Gunsalus, K. T., Wagoner, M. P., Meyer, K., Potter, W. B., Schoenike, B., Kim, S., Alexander, C. M., Friedl, A. & Roopra, A. Induction of the RNA regulator LIN28A is required for the growth and pathogenesis of RESTless breast tumors. *Cancer Res* **72**, 3207-3216, doi:10.1158/0008-5472.CAN-11-1639 (2012).
- 10 Wagoner, M. P., Gunsalus, K. T., Schoenike, B., Richardson, A. L., Friedl, A. & Roopra, A. The transcription factor REST is lost in aggressive breast cancer. *PLoS Genet* **6**, e1000979, doi:10.1371/journal.pgen.1000979 (2010).
- 11 The Cancer Genome Atlas, N., Koboldt, D. C., Fulton, R. S., McLellan, M. D., Schmidt, H., Kalicki-Veizer, J., McMichael, J. F., Fulton, L. L., Dooling, D. J., Ding, L., Mardis, E. R., Wilson, R. K., Ally, A., Balasundaram, M., Butterfield, Y. S. N., Carlsen, R., Carter, C., Chu, A., Chuah, E., Chun, H.-J. E., Cope, R. J. N., Dhalla, N., Guin, R., Hirst, C., Hirst, M., Holt, R. A., Lee, D., Li, H. I., Mayo, M., Moore, R. A., Mungall, A. J., Pleasance, E., Gordon Robertson, A., Schein, J. E., Shafiei, A., Sipahimalani, P., Slobodan, J. R., Stoll, D., Tam, A., Thiessen, N.,

Varhol, R. J., Wye, N., Zeng, T., Zhao, Y., Birol, I., Jones, S. J. M., Marra, M. A., Cherniack, A. D., Saksena, G., Onofrio, R. C., Pho, N. H., Carter, S. L., Schumacher, S. E., Tabak, B., Hernandez, B., Gentry, J., Nguyen, H., Crenshaw, A., Ardlie, K., Beroukhi, R., Winckler, W., Getz, G., Gabriel, S. B., Meyerson, M., Chin, L., Park, P. J., Kucherlapati, R., Hoadley, K. A., Todd Auman, J., Fan, C., Turman, Y. J., Shi, Y., Li, L., Topal, M. D., He, X., Chao, H.-H., Prat, A., Silva, G. O., Iglesia, M. D., Zhao, W., Usary, J., Berg, J. S., Adams, M., Booker, J., Wu, J., Gulabani, A., Bodenheimer, T., Hoyle, A. P., Simons, J. V., Soloway, M. G., Mose, L. E., Jefferys, S. R., Balu, S., Parker, J. S., Neil Hayes, D., Perou, C. M., Malik, S., Mahurkar, S., Shen, H., Weisenberger, D. J., Triche Jr, T., Lai, P. H., Bootwalla, M. S., Maglinte, D. T., Berman, B. P., Van Den Berg, D. J., Baylin, S. B., Laird, P. W., Creighton, C. J., Donehower, L. A., Getz, G., Noble, M., Voet, D., Saksena, G., Gehlenborg, N., DiCara, D., Zhang, J., Zhang, H., Wu, C.-J., Yingchun Liu, S., Lawrence, M. S., Zou, L., Sivachenko, A., Lin, P., Stojanov, P., Jing, R., Cho, J., Sinha, R., Park, R. W., Nazaire, M.-D., Robinson, J., Thorvaldsdottir, H., Mesirov, J., Park, P. J., Chin, L., Reynolds, S., Kreisberg, R. B., Bernard, B., Bressler, R., Erkkila, T., Lin, J., Thorsson, V., Zhang, W., Shmulevich, I., Ciriello, G., Weinhold, N., Schultz, N., Gao, J., Cerami, E., Gross, B., Jacobsen, A., Sinha, R., Arman Aksoy, B., Antipin, Y., Reva, B., Shen, R., Taylor, B. S., Ladanyi, M., Sander, C., Anur, P., Spellman, P. T., Lu, Y., Liu, W., Verhaak, R. R. G., Mills, G. B., Akbani, R., Zhang, N., Broom, B. M., Casasent, T. D., Wakefield, C., Unruh, A. K., Baggerly, K., Coombes, K., Weinstein, J. N., Haussler, D., Benz, C. C., Stuart, J. M., Benz, S. C., Zhu, J., Szeto, C. C., Scott, G. K., Yau, C., Paull, E. O., Carlin, D., Wong, C., Sokolov, A., Thusberg, J., Mooney, S., Ng, S., Goldstein, T. C., Ellrott, K., Grifford, M., Wilks, C., Ma, S., Craft, B., Yan, C., Hu, Y., Meerzaman, D., Gastier-Foster, J. M., Bowen, J., Ramirez, N. C., Black, A. D., Pyatt, R. E., White, P., Zmuda, E. J., Frick, J., Lichtenberg, T. M., Brookens, R., George, M. M., Gerken, M. A., Harper, H. A., Leraas, K. M., Wise, L. J., Tabler, T. R., McAllister, C., Barr, T., Hart-Kothari, M., Tarvin, K., Saller, C., Sandusky, G., Mitchell, C., Iacocca, M. V., Brown, J., Rabeno, B., Czerwinski, C., Petrelli, N., Dolzhansky, O., Abramov, M., Voronina, O., Potapova, O., Marks, J. R., Suchorska, W. M., Murawa, D., Kyrcle, W., Ibbs, M., Korski, K., Szychała, A., Murawa, P., Brzeziński, J. J., Perz, H., Łażniak, R., Teresiak, M., Tatka, H., Leporowska, E., Bogusz-Czerniewicz, M., Malicki, J., Mackiewicz, A., Wiznerowicz, M., Van Le, X., Kohl, B., Viet Tien, N., Thorp, R., Van Bang, N., Sussman, H., Duc Phu, B., Hajek, R., Phi Hung, N., Viet The Phuong, T., Quyet Thang, H., Zaki Khan, K., Penny, R., Mallery, D., Curley, E., Shelton, C., Yena, P., Ingle, J. N., Couch, F. J., Lingle, W. L., King, T. A., Maria Gonzalez-Angulo, A., Mills, G. B., Dyer, M. D., Liu, S., Meng, X., Patangan, M., Waldman, F., Stöppler, H., Kimryn Rathmell, W., Thorne, L., Huang, M., Boice, L., Hill, A., Morrison, C., Gaudioso, C., Bshara, W., Daily, K., Egea, S. C., Pegram, M. D., Gomez-Fernandez, C., Dhir, R., Bhargava, R., Brufsky, A., Shriver, C. D., Hooke, J. A., Leigh Campbell, J., Mural, R. J., Hu, H., Somiari, S., Larson, C., Deyarmin, B., Kvecher, L., Kovatich, A. J., Ellis, M. J., King, T. A., Hu, H., Couch, F. J., Mural, R. J., Stricker, T., White, K., Olopade, O., Ingle, J. N., Luo, C., Chen, Y., Marks, J. R., Waldman, F., Wiznerowicz, M., Bose, R., Chang, L.-W., Beck, A. H., Maria Gonzalez-Angulo, A., Pihl, T., Jensen, M., Sfeir, R., Kahn, A., Chu, A., Kothiyal, P., Wang, Z., Snyder, E., Pontius, J., Ayala, B., Backus, M., Walton, J., Baboud, J., Berton, D., Nicholls, M., Srinivasan, D., Raman, R., Girshik, S., Kigonya, P., Alonso, S., Sanbhadti, R., Barletta, S., Pot, D., Sheth, M., Demchok, J. A., Mills Shaw, K. R., Yang, L., Eley, G., Ferguson, M. L.,



- Tarnuzzer, R. W., Zhang, J., Dillon, L. A. L., Buetow, K., Fielding, P., Ozenberger, B. A., Guyer, M. S., Sofia, H. J. & Palchik, J. D. Comprehensive molecular portraits of human breast tumours. *Nature* **490**, 61, doi:10.1038/nature11412  
<https://www.nature.com/articles/nature11412#supplementary-information> (2012).
- 12 Carroll, J. S., Liu, X. S., Brodsky, A. S., Li, W., Meyer, C. A., Szary, A. J., Eeckhoute, J., Shao, W., Hestermann, E. V., Geistlinger, T. R., Fox, E. A., Silver, P. A. & Brown, M. Chromosome-wide mapping of estrogen receptor binding reveals long-range regulation requiring the forkhead protein FoxA1. *Cell* **122**, 33-43, doi:10.1016/j.cell.2005.05.008 (2005).
- 13 Eeckhoute, J., Keeton, E. K., Lupien, M., Krum, S. A., Carroll, J. S. & Brown, M. Positive cross-regulatory loop ties GATA-3 to estrogen receptor alpha expression in breast cancer. *Cancer Res* **67**, 6477-6483, doi:10.1158/0008-5472.CAN-07-0746 (2007).
- 14 Lacroix, M. & Leclercq, G. About GATA3, HNF3A, and XBP1, three genes co-expressed with the oestrogen receptor-alpha gene (ESR1) in breast cancer. *Mol Cell Endocrinol* **219**, 1-7, doi:10.1016/j.mce.2004.02.021 (2004).
- 15 Lin, C. Y., Vega, V. B., Thomsen, J. S., Zhang, T., Kong, S. L., Xie, M., Chiu, K. P., Lipovich, L., Barnett, D. H., Stossi, F., Yeo, A., George, J., Kuznetsov, V. A., Lee, Y. K., Charn, T. H., Palanisamy, N., Miller, L. D., Cheung, E., Katzenellenbogen, B. S., Ruan, Y., Bourque, G., Wei, C. L. & Liu, E. T. Whole-genome cartography of estrogen receptor alpha binding sites. *PLoS Genet* **3**, e87, doi:10.1371/journal.pgen.0030087 (2007).
- 16 Perou, C. M., Sorlie, T., Eisen, M. B., van de Rijn, M., Jeffrey, S. S., Rees, C. A., Pollack, J. R., Ross, D. T., Johnsen, H., Akslen, L. A., Fluge, O., Pergamenschikov, A., Williams, C., Zhu, S. X., Lonning, P. E., Borresen-Dale, A. L., Brown, P. O. & Botstein, D. Molecular portraits of human breast tumours. *Nature* **406**, 747-752, doi:10.1038/35021093 (2000).
- 17 Sorlie, T., Tibshirani, R., Parker, J., Hastie, T., Marron, J. S., Nobel, A., Deng, S., Johnsen, H., Pesich, R., Geisler, S., Demeter, J., Perou, C. M., Lonning, P. E., Brown, P. O., Borresen-Dale, A. L. & Botstein, D. Repeated observation of breast tumor subtypes in independent gene expression data sets. *Proc Natl Acad Sci U S A* **100**, 8418-8423, doi:10.1073/pnas.0932692100 (2003).
- 18 Yamaguchi, N., Ito, E., Azuma, S., Honma, R., Yanagisawa, Y., Nishikawa, A., Kawamura, M., Imai, J., Tatsuta, K., Inoue, J., Semba, K. & Watanabe, S. FoxA1 as a lineage-specific oncogene in luminal type breast cancer. *Biochem Biophys Res Commun* **365**, 711-717, doi:10.1016/j.bbrc.2007.11.064 (2008).
- 19 Zaret, K. S. & Carroll, J. S. Pioneer transcription factors: establishing competence for gene expression. *Genes Dev* **25**, 2227-2241, doi:10.1101/gad.176826.111 (2011).
- 20 Hwang, C., Giri, V. N., Wilkinson, J. C., Wright, C. W., Wilkinson, A. S., Cooney, K. A. & Duckett, C. S. EZH2 regulates the transcription of estrogen-responsive genes through association with REA, an estrogen receptor corepressor. *Breast Cancer Res Treat* **107**, 235-242, doi:10.1007/s10549-007-9542-7 (2008).
- 21 Shang, Y., Hu, X., DiRenzo, J., Lazar, M. A. & Brown, M. Cofactor dynamics and sufficiency in estrogen receptor-regulated transcription. *Cell* **103**, 843-852 (2000).
- 22 Frietze, S., O'Geen, H., Littlepage, L. E., Simion, C., Sweeney, C. A., Farnham, P. J. & Krig, S. R. Global analysis of ZNF217 chromatin occupancy in the breast cancer cell genome reveals an association with ERalpha. *BMC Genomics* **15**, 520, doi:10.1186/1471-2164-15-520 (2014).

- 23 Szklarczyk, D., Franceschini, A., Wyder, S., Forslund, K., Heller, D., Huerta-Cepas, J., Simonovic, M., Roth, A., Santos, A., Tsafou, K. P., Kuhn, M., Bork, P., Jensen, L. J. & von Mering, C. STRING v10: protein-protein interaction networks, integrated over the tree of life. *Nucleic Acids Res* **43**, D447-452, doi:10.1093/nar/gku1003 (2015).
- 24 Szklarczyk, D., Gable, A. L., Lyon, D., Junge, A., Wyder, S., Huerta-Cepas, J., Simonovic, M., Doncheva, N. T., Morris, J. H., Bork, P., Jensen, L. J. & Mering, C. V. STRING v11: protein-protein association networks with increased coverage, supporting functional discovery in genome-wide experimental datasets. *Nucleic Acids Res* **47**, D607-D613, doi:10.1093/nar/gky1131 (2019).
- 25 Szklarczyk, D., Morris, J. H., Cook, H., Kuhn, M., Wyder, S., Simonovic, M., Santos, A., Doncheva, N. T., Roth, A., Bork, P., Jensen, L. J. & von Mering, C. The STRING database in 2017: quality-controlled protein-protein association networks, made broadly accessible. *Nucleic Acids Res* **45**, D362-D368, doi:10.1093/nar/gkw937 (2017).
- 26 Fabregat, A., Jupe, S., Matthews, L., Sidiropoulos, K., Gillespie, M., Garapati, P., Haw, R., Jassal, B., K€orninger, F., May, B., Milacic, M., Roca, C. D., Rothfels, K., Sevilla, C., Shamovsky, V., Shorser, S., Varusai, T., Viteri, G., Weiser, J., Wu, G., Stein, L., Hermjakob, H. & D'Eustachio, P. The Reactome Pathway Knowledgebase. *Nucleic Acids Res* **46**, D649-D655, doi:10.1093/nar/gkx1132 (2018).
- 27 Castles, C. G., Oesterreich, S., Hansen, R. & Fuqua, S. A. Auto-regulation of the estrogen receptor promoter. *J Steroid Biochem Mol Biol* **62**, 155-163 (1997).
- 28 Fullwood, M. J., Liu, M. H., Pan, Y. F., Liu, J., Xu, H., Mohamed, Y. B., Orlov, Y. L., Velkov, S., Ho, A., Mei, P. H., Chew, E. G., Huang, P. Y., Welboren, W. J., Han, Y., Ooi, H. S., Ariyaratne, P. N., Vega, V. B., Luo, Y., Tan, P. Y., Choy, P. Y., Wansa, K. D., Zhao, B., Lim, K. S., Leow, S. C., Yow, J. S., Joseph, R., Li, H., Desai, K. V., Thomsen, J. S., Lee, Y. K., Karuturi, R. K., Herve, T., Bourque, G., Stunnenberg, H. G., Ruan, X., Cacheux-Rataboul, V., Sung, W. K., Liu, E. T., Wei, C. L., Cheung, E. & Ruan, Y. An oestrogen-receptor-alpha-bound human chromatin interactome. *Nature* **462**, 58-64, doi:10.1038/nature08497 (2009).
- 29 Carroll, J. S., Meyer, C. A., Song, J., Li, W., Geistlinger, T. R., Eeckhoutte, J., Brodsky, A. S., Keeton, E. K., Fertuck, K. C., Hall, G. F., Wang, Q., Bekiranov, S., Sementchenko, V., Fox, E. A., Silver, P. A., Gingeras, T. R., Liu, X. S. & Brown, M. Genome-wide analysis of estrogen receptor binding sites. *Nat Genet* **38**, 1289-1297, doi:10.1038/ng1901 (2006).
- 30 Sams, D. S., Nardone, S., Getselter, D., Raz, D., Tal, M., Rayi, P. R., Kaphzan, H., Hakim, O. & Elliott, E. Neuronal CTCF Is Necessary for Basal and Experience-Dependent Gene Regulation, Memory Formation, and Genomic Structure of BDNF and Arc. *Cell Rep* **17**, 2418-2430, doi:10.1016/j.celrep.2016.11.004 (2016).
- 31 Sawado, T., Igarashi, K. & Groudine, M. Activation of beta-major globin gene transcription is associated with recruitment of NF-E2 to the beta-globin LCR and gene promoter. *Proc Natl Acad Sci U S A* **98**, 10226-10231, doi:10.1073/pnas.181344198 (2001).
- 32 Andres, M. E., Burger, C., Peral-Rubio, M. J., Battaglioli, E., Anderson, M. E., Grimes, J., Dallman, J., Ballas, N. & Mandel, G. CoREST: a functional corepressor required for regulation of neural-specific gene expression. *Proc Natl Acad Sci U S A* **96**, 9873-9878 (1999).
- 33 Garriga-Canut, M., Schoenike, B., Qazi, R., Bergendahl, K., Daley, T. J., Pfender, R. M., Morrison, J. F., Ockuly, J., Stafstrom, C., Sutula, T. & Roopra, A. 2-Deoxy-D-glucose reduces

- epilepsy progression by NRSF-CtBP-dependent metabolic regulation of chromatin structure. *Nat Neurosci* **9**, 1382-1387, doi:10.1038/nn1791 (2006).
- 34 Huang, Y., Myers, S. J. & Dingledine, R. Transcriptional repression by REST: recruitment of Sin3A and histone deacetylase to neuronal genes. *Nat Neurosci* **2**, 867-872, doi:10.1038/13165 (1999).
- 35 Roopra, A., Sharling, L., Wood, I. C., Briggs, T., Bachfischer, U., Paquette, A. J. & Buckley, N. J. Transcriptional repression by neuron-restrictive silencer factor is mediated via the Sin3-histone deacetylase complex. *Mol Cell Biol* **20**, 2147-2157 (2000).
- 36 Dietrich, N., Lerdrup, M., Landt, E., Agrawal-Singh, S., Bak, M., Tommerup, N., Rappsilber, J., Sodersten, E. & Hansen, K. REST-mediated recruitment of polycomb repressor complexes in mammalian cells. *PLoS Genet* **8**, e1002494, doi:10.1371/journal.pgen.1002494 (2012).
- 37 McGann, J. C., Oyer, J. A., Garg, S., Yao, H., Liu, J., Feng, X., Liao, L., Yates, J. R., 3rd & Mandel, G. Polycomb- and REST-associated histone deacetylases are independent pathways toward a mature neuronal phenotype. *Elife* **3**, e04235, doi:10.7554/eLife.04235 (2014).
- 38 Mourad, R. & Cuvier, O. Computational Identification of Genomic Features That Influence 3D Chromatin Domain Formation. *PLoS Comput Biol* **12**, e1004908, doi:10.1371/journal.pcbi.1004908 (2016).