

# Network Dynamics of ER- $\alpha$ activation reveals reprogramming of GRHL2

Andrew N. Holding<sup>1,3,\*</sup>, Federico M. Giorgi<sup>1,3</sup>, Amanda Donnelly<sup>1</sup>, Amy E. Cullen<sup>1</sup>, Luke Selth<sup>2</sup>,  
and Florian Markowitz<sup>1</sup>

<sup>1</sup> CRUK Cambridge Institute, University of Cambridge, Robinson Way, Cambridge, CB2 0RE

<sup>2</sup> Dame Roma Mitchell Cancer Research Laboratories and Freemasons Foundation Centre for Men's Health, Adelaide Medical School, The University of Adelaide, SA, Australia

<sup>3</sup> These authors contributed equally to this work.

\*Corresponding Author

## Abstract

**Background** Estrogen Receptor-alpha (ER) is the main driver of ~75% of all breast cancers. Upon stimulation, ER forms a complex on the chromatin at enhancers and promoters that leads to increased transcription of nearby genes. A critical feature of ER action is a cyclical binding pattern with a period of 90 minutes. However, analysis of ER binding dynamics has so far been restricted to the promoters of individual target genes. It is unknown how cyclical ER binding occurs genome-wide and whether this phenomenon is influenced by ER cofactors.

**Results** Here, we present a novel approach to dissect the regulation of ER activation based on network analysis of time-course genome-wide DNA binding data. We measured ER binding by ChIP-Seq at three timepoints (0, 45, 90 minutes) and developed an approach, called VULCAN, to overlay this binding information on a gene coexpression network. We benchmarked our approach in a comparison study and confirmed that it rediscovered known components of the ER signalling axis. Using VULCAN, we found that the activation ER results in the reprogramming of the transcription factor GRHL2 and independently validated this result by ChIP-seq and quantitative proteomics (qPLEX-RIME). Further, E2-responsive GRHL2 binding was found to be concurrent with an ER-responsive increase in eRNA transcription, and we show GRHL2 negatively regulates transcriptional activity at these sites.

**Conclusions** We present a general framework to predict key regulatory proteins from differential transcription factor binding data, which uncovered that activation of the ER leads to reprogramming of GRHL2.

**Keywords:** Breast Cancer, Network Analysis, Dynamics, ER, Master Regulator, ChIP-Seq, VULCAN, GRHL2, Squelching, P300

# Background

Breast cancer is the most common form of cancer in women in North America and Europe. The majority of breast cancers are associated with the deregulation of Estrogen Receptor-alpha (ER), which drives the growth and proliferation of the tumor. ER is the key prognostic marker in the clinic and the target of the first lines of treatment in estrogen receptor positive (ER+) tumors. ER-targeting pharmacological strategies include SERMs (selective estrogen receptor modulators) and SERDs (estrogen receptor degraders) e.g. Tamoxifen and Fulvestrant, or aromatase inhibitors that block the production of estrogens in the body [1].

## **The role of ER has been extensively studied genome-wide**

On activation, ER binds to promoter and enhancer regions containing Estrogen-Response Elements (ERE) [2] to stimulate the transcription of hundreds of genes [3,4]. Gene expression is driven by both the recruitment of the basal transcription machinery to these loci and through longer range interactions [5]. Analysis of ER-target genes showed that many are proliferative in function and drive the growth of the tumor [6]

## **ER associates with a wide range of cofactors**

On the treatment of ER+ cells with estra-2-diol (E2), ER recruits several cofactors to form a complex on the chromatin. FOXA1 is of particular interest as the protein shares nearly 50% of its genomic binding sites with ER and has been shown to operate as a pioneer factor before ER activation [7]. It is through FOXA1 and other cofactors [8,9], e.g. SRC-1, that ER is able to recruit RNA Polymerase II at the gene promoter sites in order to initiate transcription [10].

## The binding of ER to chromatin is highly dynamic

Early studies have shown that after stimulation with E2, ER binding to EREs can occur within minutes [11]. Maximum ER occupancy at promoters of target genes such as CTSD and TFF1 is achieved after 45' in MCF7 cells. Roughly 90' after estradiol treatment, ER is partially released from the promoters of the target gene (**Figure 1A**), only to reoccupy the site again at 135' and release them at 180', in a 45' phase cyclical manner [12]. The 90' occupancy phase has been shown to be independent of new protein synthesis, and is therefore thought to operate at the post-translational level. This cyclic, proteasome-mediated turnover of unliganded and liganded ER $\alpha$  on responsive promoters is an integral feature of estrogen signaling [13,14].

## Network analysis to infer TF activity

Given the three features discussed above, ER is a prime target for systems biology. Usage of gene regulatory networks to analyze biological systems has witnessed an exponential increase in the last decade, due to the ease of obtaining genome-wide expression data [15–17]. Recently, the VIPER approach to interrogate these network models has been proposed to infer transcription factor activity using the expression of a collection of their putative targets, *i.e.* their regulon [18]. In the VIPER algorithm, gene-level *differential expression* signatures are obtained for either individual samples (relative to the mean of the dataset) or between groups of samples, and regulons are tested for enrichment.

In our study, we propose an extension of the VIPER algorithm to specifically analyze TF occupancy in ChIP-Seq experiments. Our algorithm, called “**VirtUaL ChIP-Seq Analysis through Networks**” (VULCAN), uses ChIP-Seq data obtained for a given TF to provide candidate coregulators of the response to a given stimulus (**Figure 2**). The analysis is based on identifying differentially bound genes and testing their enrichment in the regulon of potential co-regulatory

factors. By applying VULCAN to ChIP-seq time-course data of ER activation, our study provides new temporal insights into ER cofactors on a genome-wide scale.

## Results

The aim of our study is to identify key coregulators of the ER binding process and led to the development of VULCAN, that infers changes in cofactor activity from differential ChIP-seq analysis using network analysis. An implementation of VULCAN in R is available on Bioconductor.org [<https://bioconductor.org/packages/release/bioc/html/vulcan.html>] and the scripts to replicate our analysis are available as a supplementary Rmarkdown file. Unless otherwise specified, all p-values were Bonferroni-corrected.

## Network Inference

We generated co-expression networks using the most recent implementation of ARACNe [19] on the METABRIC dataset [20] and the independent TCGA data set [21]. Briefly, ARACNe generates gene networks by estimating putative regulatory interactions between transcription factors and target genes using mutual information between gene expression profiles. As an example, **Figure 2C** shows an ARACNe-inferred targets of ESR1. The sets of targets of each TF, its regulons, were merged in a genome-wide transcriptional regulation network, as shown in a minimal diagram in **Figure 2D**.

### **Regulatory network analysis to detect ER cofactors**

In order to understand which co-factors could be responsible for the temporal behavior of ER, we modified the VIPER algorithm [18] to perform master regulator analysis with differential

binding signatures. Master regulator analysis [22] is an algorithm developed to identify transcription factors whose regulon is enriched within a list of differentially expressed genes. VULCAN extends the VIPER approach to use differential binding profiles rather than differential expression profiles. In this way, VULCAN can test the enrichment of TF regulons in ER occupancy signatures derived from ChIP-Seq experiments (**Figure 2E**).

## Benchmarking VULCAN

### Comparing VULCAN's Mutual Information and Partial Correlation networks

VULCAN uses mutual information networks like VIPER [18]. To test the robustness of our approach to different underlying networks, we compared mutual information networks with partial correlation networks using different correlation thresholds. We generated several partial correlation networks from the TCGA data using the same input as the ARACNe network used by VULCAN. We tested the overlap of every partial correlation network with the ARACNe network using the Jaccard Index (JI) criterion (**Suppl. Figure 33**). Finally, we show how the Jaccard Index between partial correlation networks and the ARACNe network is always significantly higher than expected by selecting random network edges (**Suppl. Figure 34**). This confirms previous observations that partial correlation and mutual information networks are highly similar [23]

### Comparing VULCAN with alternative methods for target enrichment analysis

We compared VULCAN's GSEA approach with three independent methods previously applied to benchmark VIPER [18]. The first implemented a t-test based method, which takes the targets of a TF and integrates their p-value in a specific contrast. The method is similar to VIPER but involves a Fisher p-value integration step. The integrated test lacks stringency and results in

nearly all regulons as significantly enriched (**Suppl. Figure 35**). Second, we implemented a fraction of targets method, defining for every TF the fraction of their targets that are also differentially bound. This alternative to VULCAN ignores the MI strength of interaction and the individual strengths of differential bindings, reducing the resolving power of the algorithm (**Suppl. Figure 36**). Finally, we compared to a Fisher's Exact Method which assesses the overlap between networks and significant differential binding. This method is too stringent (as observed in the original VIPER paper) [18] and even without p-value correction there are no significant results, even at low stringency. In short, our analysis demonstrates the low sensitivity of this method (**Suppl. Figure 37**). In summary all three alternative methods we tested (t-test based; fraction of targets method; and Fisher's Exact Method) all resulted in reduced performance on our dataset compared to VULCAN.

### **Comparing VULCAN with Online Tools (GREAT, ISMARA & CHIP-Enrich)**

To further validate our method, we compared the output of our GSEA analysis with different versions of promoter-enrichment approaches implemented by GREAT [24], ISMARA [25] and CHIP-Enrich [26]. The VULCAN analysis shows a significant overlap in terms of significant pathways with the GREAT method (**Suppl. Figure 38**). CHIP-enrich computes enrichment for a number of TFs which are amongst the most significant in VULCAN, but it fails at identifying ESR1 as the top Transcription Factor affected by our experiment (**Suppl. Figure 39**). ISMARA succeeds at identifying ESR1 using a motif-based analysis, but does not identify other candidate binding TFs, as expected, being the experiment targeted at the estrogen receptor (**Suppl. Figure 40**). In summary, VULCAN outperforms both CHIP-Enrich and ISMARA. In-terms of pathway analysis VULCAN reassuringly provides significantly overlaps with GREAT while our network analysis contributes additional value through inference of TF factor activity.

## VULCAN analysis of ER activation

### Differential Binding Analysis

We performed four replicated ChIP-Seq experiments for ER at three timepoints: 0, 45 and 90 minutes after estradiol treatment (**Figure 1**). The binding profile of ER at each time point was then compared between timepoints using differential binding analysis.

Differential Binding Analysis (**Figure 1B,C**) identified 18,900 statistically significant binding events at 45 minutes ( $p < 0.05$ ). We observed the previously reported reduction of ER binding at 90 minutes [13] on a genome-wide level (17,896 significant binding events), but with a smaller amplitude than previous gene-specific assays [14].

We performed motif enrichment analysis (HOMER software) on ER binding sites detected by differential binding analysis. This analysis confirmed a strong enrichment for a single element, ERE, bound at both 45 and 90 minutes, with a corrected p-value of 0.0029 (**Figure 3F**). When clustered according to peak intensity, samples cluster tightly in two groups: treated and untreated (**Suppl. Figures 2, 3 and 4**), but treatment at 45 and 90 minutes is detectably different on a genome-wide scale, as highlighted by Principal Component Analysis (**Suppl. Figures 5 and 6**).

A potential cause for the smaller amplitude we observed in ER cycling may relate to the fact that ChIP-Seq is not inherently quantitative, and hence the typical normalisation strategies applied to ChIP-Seq data are likely to suppress global changes [27,28]. We therefore validated the ER binding behavior with ChIP-qPCR (**Figure 1A**) and observed the same reduction in amplitude at specific binding events as was predicted by ChIP-seq. Another potential reason for the difference in amplitude was we did not treat  $\alpha$ -Amanitin prior to treatment with E2 [14] as this perturbation would further separate the experimental condition from clinical interpretation.



### **VULCAN groups genes by temporal dynamics of ER binding**

We leveraged the information contained in mutual information networks to establish TF networks enriched in the differential binding patterns induced by estradiol. From the temporal comparison ER binding we established four classes of binding pattern: early responders, repressed transcription factors, late responders and candidate cyclic genes (**Figure 3**).

Using VULCAN, we defined TF network activity of occupied regulatory regions (**Figure 3A**) according to the binding of the ER within their promoter and enhancer regions (10kb upstream of the Transcription Starting Site). We define as early responders TFs whose network is upregulated at both 45 and 90 minutes (**Figure 3B**): these genes include AR, SP1 and CITED1. TFs with opposite behavior (namely, TFs whose negative/repressed targets in the ARACNe model are occupied by ER), or “repressed TFs” include GLI4, MYCN and RAD21 (**Figure 3C**). Some TFs appear to have their targets bound at 45 minutes, but then unoccupied at 90 minutes. This “updown” behavior is consistent with the cyclic properties of certain components of the ER DNA-binding complex observed previously, and therefore we dubbed them “candidate cyclic TFs” (**Figure 3D**). We also define a “late responders” category, expecting to find TFs active at 90 minutes but not at 45 minutes. While this category exists, it is just below the significance threshold at 45 minutes, and notably it contains both ESR1 and the known ESR1 interactor GATA3 (**Figure 3E**).

### **Validating VULCAN results on independent data**

We repeated the analysis in Figure 3A-E on independent data from TCGA (**Suppl. Figures 16–20**). We again found enrichment for ESR1, again as a later responder. GATA3 was detected as a early responder, likely as a result of the METABRIC result being only slightly below significance threshold at 45 minutes. Notably PGR was shifted from an early responder to a late responder.

To ensure the robustness of results we performed a joint analysis of data obtained from both networks. At 45 minutes (**Figure 4A**) and 90 minutes (**Figure 4B**) we identified robust candidates: specifically, ESR1, GATA3 and RARA networks, amongst others, were consistently occupied by ER in both time points. On the other hand, some genes, including HSF1 and GRHL2, were significantly repressed in the joint analysis.

As a negative control, we used a different context ARACNe network, derived from the TCGA AML dataset. This network shows globally weaker enrichment scores and a weak positive correlation with the results obtained through breast cancer regulatory models (**Suppl. Figure 22**).

### **Pathway analysis of regulatory region binding**

We performed a Gene Set Enrichment Analysis [29] and an associated Rank Enrichment Analysis [18] using the differential binding at gene regulatory regions (with time 0 as reference). Individual differential binding signatures for GSEA were calculated using a negative binomial test implemented by DiffBind [30]. The collective contribution of differentially bound sites highlights several pathways ER-related pathways [31–33] (**Suppl. Figure 23**) in both the GSEA and aREA analyses. The strongest pathway upregulated pathway in both time points (**Table S2 and S5**) was derived via RNA-Seq in an MCF7 study using estradiol treatment [32] confirming the reproducibility of our data set.

## Validating VULCAN results by quantitative proteomics

We tested the performance of VULCAN against a complementary experimental approach called RIME [34] combined with TMT [35] (qPLEX-RIME), which aims at identifying interactors of ER within the ER-chromatin complex. We generated ER qPLEX-RIME data from MCF7 cells treated with estradiol at both 45 and 90 minutes and compared this with the VULCAN dataset, with the aim of identifying TFs upstream of the observed differential binding (**Suppl. Figure 31**). We found known ESR1 interactors with both methods, namely HDAC1, NCOA3, GATA3 and RARA with positive Network Enrichment Score (NES) [18] and GRHL2 with a negative NES.

## The GRHL2 Transcription Factor

In our analysis of ER dynamics the GRHL2 transcription factor stood out. In both the METABRIC and TCGA networks GRHL2 was significantly repressed, yet in our proteomics analysis the protein was significantly increased. Therefore we set out to validate experimentally GRHL2 as an ESR1 cofactor, possibly with repression properties for the ER complex

Our analysis shows that the genes occupied by the ER complex do not form part of the GRHL2 regulon, using both TCGA-derived and METABRIC-derived regulatory models. Our analysis highlights the small overlap between the ESR1 (Estrogen Receptor) and GRHL2 networks (**Suppl. Figure 28**), hinting at complementary signals not dependent on global network overlaps. In fact, there is merely a weak, positive correlation between ESR1 and GRHL2 expression in the TCGA breast cancer dataset (**Suppl. Figure 29**) and also in the METABRIC breast cancer dataset (**Suppl. Figure 30**). Furthermore, GRHL2 has a visibly lower variance than ESR1: it does not change significantly in different PAM50 subtypes, although it is lower in normal compared to malignant tissue.

The low overlap between networks and the low correlation in expression profiles could be explained by the fact that the GRHL2 role may not be carried out by its transcription, and rather being controlled by other mechanisms like phosphorylation, subcellular localization or on-chromatin interactions.

### **GRHL2 activity before and after stimulation with E2**

Differential ChIP-seq analysis of GRHL2 binding between 0 and 45 minutes indicated that GRHL2 binding is altered on treatment with E2 (**Figure 5A**). VULCAN analysis of the GRHL2 differential binding showed a consistent Network Enrichment Score for both the TCGA- and METABRIC-derived networks for ER, but not for FOXA1 or GRHL2 (**Figure 5B**). Individual analysis of peaks show that typically ER promoter sites, e.g. RARa, were not the target of this redistribution of GRHL2, as these sites were occupied by GRHL2 before E2 stimulation. We propose that GRHL2's occupancy at these site is via a direct binding at FOXA1 sites as previously described [36]. Motif analysis of the sites within increased GRHL2 occupancy showed enrichment for the full ERE (p-value =  $1 \times 10^{-86}$ ) and the GRHL2 binding motif (p-value =  $1 \times 10^{-31}$ ).

qPLEX-RIME analysis of GRHL2 interactions showed high coverage of the bait protein (>59%) and, in both the estrogen-free and estrogenic conditions, high levels of transcription-related protein interactors including HDAC1 (p-value =  $6.4 \times 10^{-9}$ ), TIF1A (p-value =  $6.4 \times 10^{-9}$ ), PRMT (p-value =  $6.4 \times 10^{-9}$ ) and GTF3C2 (p-value =  $4.6 \times 10^{-9}$ ). P-values given for estrogenic conditions, estrogen-free conditions were comparable. Comparing the GRHL2 interactome between estrogen-free and estrogenic conditions only identified ER as a differentially bound protein that was also enriched over IgG control. Activation of ER, therefore, does not alter the majority of GRHL2 protein interactions.

The specific nature of the detected ER-GRHL2 interaction suggests that ER does not regulate or interact with GRHL2 through the recruitment of alternative cofactors. This implies instead that ER recruits GRHL2 to certain genomic loci for a specific function. We therefore undertook a comparison of GRHL2 binding with public data sets. **(Figure 5C)**. Our analysis showed that GRHL2 sites that responded to E2 were enriched for ER binding sites (in agreement with our qPLEX-RIME data) and FOXA1 (compatible with either an ER interaction or the previously reported interaction with MLL3 [36]). To establish therefore if the reprogramming of GRHL2 was primarily related to a transcriptional function or the previously described interaction with MLL3, we overlapped our GRHL2 data with that of published H3K4me1/3 [36] and P300 [37] cistromes. While H3K4me occupancy was consistent between conditions, we found P300 binding to be enriched at the E2 responsive GRHL2 sites. A more detailed analysis of the GRHL2 overlap with P300 sites showed the greatest co-occupancy of GRHL2/P300 sites was when both TFs were stimulated by E2 **(Figure 5D)**. Moreover, overlap of GRHL2 peaks with ER ChIA-PET data [ENCSR000BZZ] showed that the GRHL2 responsive sites were enriched at enhancers over promoters **(Figure 5E)**.

These findings suggested that the GRHL2-ER interaction was involved in transcription at ER enhancer sites. To explore this concept further, we investigated the transcription of enhancer RNAs at these sites using publicly available GRO-seq data [38] [GSE43836] **(Figure 5F)**. At E2 responsive sites, eRNA transcription was strongly increased by E2 stimulation; by contrast, eRNA transcription was largely independent of E2 stimulation when the entire GRHL2 cistrome was considered. Analysis of a second data set, GSE45822, corroborates these results **(Suppl. Figure S92)**.

Analysis of eRNA expression at the GREB1, TFF1 and XBP1 enhancers after over-expression of GRHL2 showed a visible decrease in eRNA production **(Figure 6)**. The

reduction in eRNA levels was significant at both the TFF1 and XBP1 enhancers ( $p < 0.05$ , pair-sample wilcoxon test). Likewise, eRNA production after GRHL2 knockdown showed a moderate increase in eRNA levels at the TFF1, XBP1 and GREB1 enhancers (**Suppl. Figure 106**). Combining data from all three sites results established the effect as significant ( $p = 0.04$ , one-tailed paired-sample wilcoxon test). Collectively, these data demonstrate that GRHL2 constraints specific ER enhancers.

## Discussion

### VirtUaL ChIP-Seq Analysis through Networks – VULCAN

Our study aimed to address the question of how ER activation is regulated at a genome-wide level. To achieve this aim, we developed and applied a network algorithm (VULCAN) that uses ChIP-seq data to reliably predict key regulatory transcription factors that impact on a TF of interest. The VULCAN algorithm was tested over a three time-point ChIP-Seq dataset of cell lines treated with E2 to test short-term dynamics of ER binding. Our analysis allowed us to provide a functional categorization of genes regulated by ER, assigning them to four different time-dependent dynamical response groups. VULCAN provided a list of TFs most likely responsible for the time-dependent responses. Supporting the robustness of the approach, VULCAN identified ER targets as the strongest responders, together with other known ER complex cofactors.

The VULCAN algorithm can be applied generally to ChIP-Seq for the identification of new key regulator interactions. Our method provides a novel approach to investigate chromatin

occupancy of cofactors that are too transient or for which no reliable antibody is available for direct ChIP-Seq analysis.

## Reprogramming of GRHL2 by ER

In the 4T1 tumor model GRHL2 was found to be significantly down-regulated in cells that had undergone EMT [39]. The same study showed that knockdown of GRHL2 in MCF10A – an ER-negative cell line– lead to loss of epithelial morphology. Overall, this suggested that the GRHL2 transcription factor plays an essential role in maintaining the epithelial phenotype of breast cells. Similar results were observed with the MDA-MB-231 model, where expression of GRHL2 resulted in reversal of EMT [40]. This result has been recapitulated in hepatocytes, where GRHL2 was found to suppress EMT by inhibiting P300 [41]. Combined these demonstrate a significant role for GRHL2 in the progression of breast cancer.

Survival Data for ER+ breast cancer (KMplotter, use gene expression,  $p=0.001$ ) and ER- (KMplotter, use gene expression,  $p=0.035$ ) shows that high GRHL2 has a negative impact of survival time in both contexts. The ability to suppress EMT has also been noted in prostate cancer, another cancer driven by a steroid hormone receptor (AR), and the genes regulated by GRHL2 are linked to disease progression [42]

Analysis of clinical samples for breast cancer tumours compared with normal cells show that ER binding location is enough to categorise the samples, and that binding sites in tumour samples were enriched with the GRHL2 motif [43].

In breast cancer, GRHL2 has previously been shown to directly interact with FOXA1, which may contribute to tethering of the histone methyltransferase MLL3 and, consequently, epigenetic marks at GRHL2/FOXA1 binding sites [36]. Our analysis, however, showed no particular enrichment for H3K4me1/3 marks at E2 responsive GRHL2 sites compared to other

GRHL2 binding sites and our proteomic analysis of interactors showed a strong association with proteins related to transcription. Therefore, while GRHL2 ChIP-seq analysis shows that GRHL2 is already bound to a proportion of FOXA1 sites before treatment of cells with E2, we still saw an increase in GRHL2 binding at 45 minutes. We proposed that these ER responsive sites are related to a role of GRHL2 in a transcriptional process independent of its interaction with MLL3. This was supported by evidence of a significant overlap with binding of the coactivator P300 and a pronounced increase in eRNA transcription on activation at E2 responsive GRHL2 sites.

Further, over-expression of GRHL2 resulted in a significant decrease in eRNA production at TFF1 and XBP1 enhancer sites. These results are consistent with previous findings that GRHL2 inhibits P300 [41] and, while the ER complex results in the activation of eRNA transcription at these sites, that GRHL2 plays a role to in fine-tuning or modulating this process.

## Conclusions

VULCAN is built on state-of-the-art network analysis tools previously applied to RNA-Seq data. By adapting network-based strategies to ChIP-Seq data, we have been able to reveal novel information regarding the regulation of breast cancer in a model system.

The VULCAN method is valuable for the discovery of transcription factors which have a role in the regulation of protein complexes that would otherwise remain hidden. The challenge of highlighting cofactors from a ChIP-Seq experiment lays in the infeasibility of reliable proteomic characterization of DNA-bound complexes at specific regions. On the other hand, while RNA-Seq is arguably the most efficient technique to obtain genome-wide quantitative measurements, any transcriptomic approach cannot provide a full picture of cellular responses



for stimuli that are provided on a shorter timescale than mRNA synthesis speed, such as the estradiol administration described in our study. VULCAN, by combining RNA-Seq derived networks and ChIP-Seq cistrome data, aims at overcoming limitations of both. Most notably, our method can work in scenarios where candidate cofactors do not have a well characterized binding site or do not even bind DNA directly.

By developing VULCAN, we have been able to rediscover known cofactors of the estradiol-responsive ER complex and predict and experimentally validate a novel protein-protein interaction.

VULCAN enabled us to identify the reprogramming of GRHL2 by the ER on stimulation with E2. Further analysis showed the process to be unrelated to the previously reported interaction with FOXA1 and MLL3 [36]. Our conclusion was that GRHL2 has a second, previously undescribed, role: negatively regulating levels of transcription at estrogen responsive enhancers (**Figure 7**). Given the central role of the ER in breast cancer development and GRHL2's own ability to regulate EMT, the discovery that ER recruits GRHL2, leading to the constraint of eRNA transcription, is an important step in enhancing our understanding of breast cancer and tumorigenesis.

## Methods

### Sample preparation

MCF7 cells were obtained from the CRUK Cambridge Institute collection, authenticated by STR genotyping and confirmed free of mycoplasma. All cells were maintained at 37 °C, 5% CO<sub>2</sub>. For each individual ChIP pull-down, we cultured 8 x 10<sup>7</sup> MCF7 cells (ATCC) across four 15 cm

diameter plates in DMEM with 10% FBS, Glutamine and Penicillin/Streptomycin (Gibco). Five days before the experiment, the cells were washed with phosphate buffered saline (PBS) and the media was replaced with clear DMEM supplemented with charcoal treated serum. The media was refreshed every 24 hours, which halted the growth of the cells and ensured that majority ER within the cell was not active. On day 5, the cells were treated with estradiol (100 nM). At the appropriate time point, the cells were washed with ice cold PBS twice and then fixed by incubating with 10mL per plate of 1% formaldehyde in unsupplemented clear media for 10 minutes. The reaction was stopped by the addition of 1.5mL of 2.5 M glycine and the plates were washed twice with ice cold PBS. Each plate was then scraped in 1 mL of PBS with protease inhibitors (PI) into a 1.5 mL microcentrifuge tube. The cells were centrifuged at 8000 rpm for 3 minutes at 4 °C and the supernatant removed. The process was repeated for a second wash in 1 mL PBS+PI and the PBS removed before storing at -80 °C.

## ChIP-Seq

Frozen samples were processed using established ChIP protocols [44] to obtain DNA fragments of ~300 bp in length. The libraries were prepared from the purified DNA using a ThruPLEX DNA-seq kit (Rubicon Genomics) and sequenced on the Illumina HiSeq Platform. Sequencing data is available from Gene Expression Omnibus, accession GSE109820.

## Differential binding analysis

Sequencing data was aligned using BWA[45] to the human genome (hg19). Reads from within the DAC Blacklisted Regions was removed before peak calling with MACS 2.1 [46] on default

parameters. The aligned reads and associated peak files were then analyzed using DiffBind [30] to identify significant changes in ER binding.

## Gene Set Enrichment Analysis (GSEA)

Gene Set Enrichment Analysis (GSEA) was performed as described by *Subramanian et al.* [47] using the curated pathway collection (C2) from MSIGDB v 5.0 with 1000 set permutations for each pathway investigated, followed by Benjamini Hochberg P-value correction.

## Motif analysis

Motif analysis of the binding regions was undertaken with Homer v4.4 [48] using default parameters. Motif logo rendering was performed using Weblogo v2.8.2 [49]

## VULCAN analysis

We reconstructed a regulatory gene network using ARACNe-AP as described by Alvarez [23]. RNA-Seq breast cancer data was downloaded from TCGA on January 2015 and VST-Normalized as described by Anders and Huber [50]. The ARACNe transcriptional regulation network was imported into R using the *viper* BioConductor package and it was interrogated using the differential binding profiles from our CHIP-Seq experiment as signatures, 45' vs control and 90' vs control. The peak-to-promoter assignment was performed using a 10kb window with respect to the transcription starting site (TSS) of every gene on the hg19 human genome. The algorithm *msVIPER* (multi-sample Virtual Inference of Protein activity by Enriched Regulon analysis) was then applied, leveraging the full set of eight replicates per group, with 1000 signature permutations and default parameters.

## qPLEX-RIME

Samples were prepared as previously described for RIME[34], protocol was modified to include TMT isobaric labels for quantification (manuscript under review). qPLEX-RIME data and analysis pipeline is available as part of the supplementary Rmarkdown.

## TF Binding Overlap

Publically available data was downloaded as described in the source publication [36–38,51] and overlap was calculated with bedtools (v2.25.0). Presented data was normalised as a percentage of GRHL2 sites.

## eRNA quantification

MCF7 cells were transfected with Smart Pool siRNA (Dharmacon, L-014515-02), siControl, GRHL2 overexpression vector (Origene, RC214498) or empty control vector using Lipofectamine 3000 (Thermo Fisher Scientific) according to the manufacturer's protocol in 6 well format. Expression was monitored by rtPCR using TaqMan assay with GAPHD as a control transcript. Knockdown efficiency was ~75% and the GRHL2 over-expression vector led a 730-fold increase in expression over control plasmid. 1 ug of purified RNA was reverse transcribed with Superscript III reverse transcriptase (Thermo Fisher Scientific, 18080085) using random primers (Promega, C1181) according to manufacturer instructions. eRNAs were

quantified with qPCR using Power SYBR™ Green PCR Master Mix (Thermo Fisher Scientific, 4367660) and denoted as relative eRNA levels after normalizing with UBC mRNA levels.

<b>Primer name</b>	<b>Sequences</b>	<b>Reference</b>
eGREB1 F	ACTGCGGCATTTCTGTGAGA	This study
eGREB1 R	ACTGCAGTTTGCCTGTCACT	This study
eXBP1 F	TGTGAGCACTTGGCATCCAT	Nagarajan et al 2014
eXBP1 R	ACAGGGCCTCATTCTCCTCT	Nagarajan et al 2014
eTFF1 F	AGGGGATGTGTGTGAGAAGG	Li et al 2013
eTFF1 R	GCTTCGAGACAGTGGGAGTC	Li et al 2013
UBC F	ATTTGGGTCGCGGTTCTTG	Peña et al 2009
UBC R	TGCCTTGACATTCTCGATGGT	Peña et al 2009

## Abbreviations

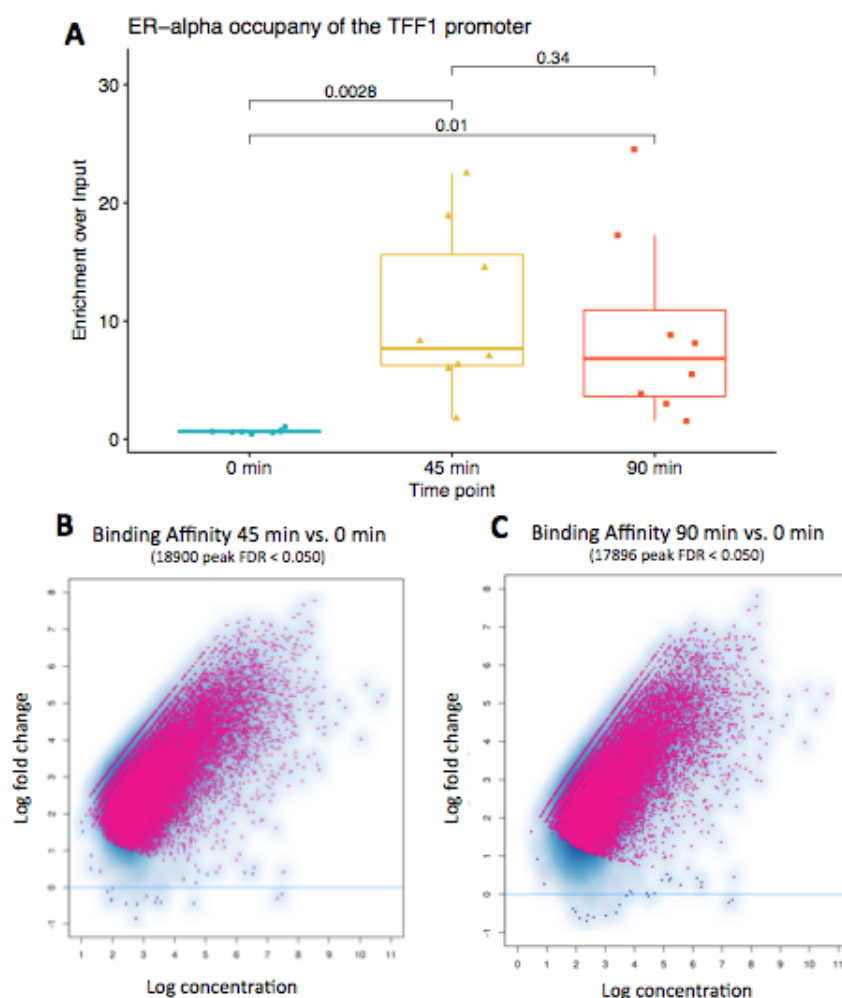
AR	Androgen Receptor
ARACNe-AP	AccuRate Algorithm for reConstruction of Network through Adaptive Partitioning
CCLE	Cancer Cell Line Encyclopedia
ChIP	Chromatin ImmunoPrecipitation
ER	Estrogen Receptor-Alpha
ERE	Estrogen-Response Elements
GSEA	Gene Set Enrichment Analysis
GRHL2	Grainyhead Like Transcription Factor 2
METABRIC	MoLEcular TAXonomy of BReast cancer International Consortium
PBS	Phosphate Buffered Saline
PI	Protease Inhibitors
PR	Progesterone Receptor
TCGA	The Cancer Genome Atlas
TF	Transcription Factor
VIPER	Virtual Inference of Protein activity by Enriched Regulon analysis
VULCAN	VirtUaL ChIP-Seq Analysis through Networks

## Declarations

Parts of this work were funded by CRUK core grant C14303/A17197 and A19274 (to FM), and by Breast Cancer Now Award (2012NovPR042).

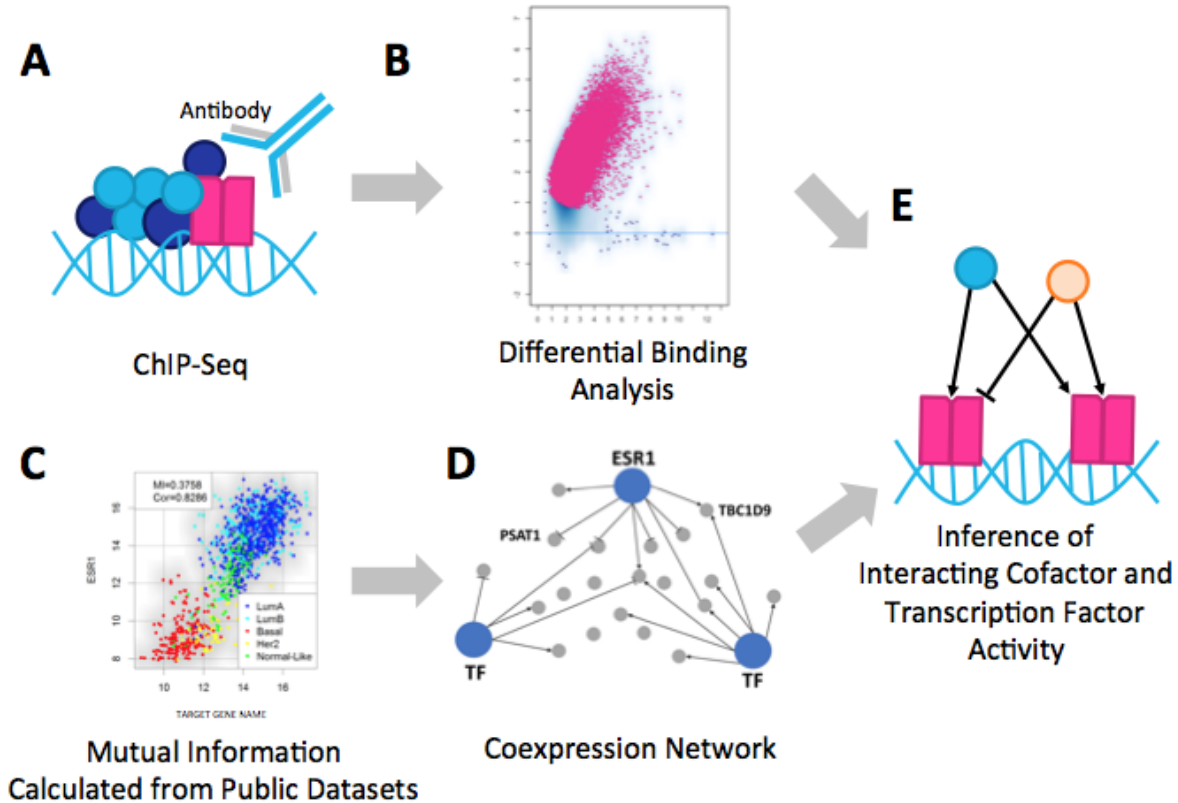
## Figure Legends

Figure 1: Dynamic behaviour during early activation of ER



ChIP-qPCR of the TFF1 gene (A) at 3 time points shows increased binding of ER at 45 minutes after MCF7 cells are stimulated by estradiol. The previously reported maximum is followed by a decrease in the TFF1 promoter occupancy at 90 minutes. P-values are generated by one-tailed t-test. The maximal point at 90 minutes was identified as an outlier ( $> \text{median} + 2 \times \text{IQR}$ ); however removal did not alter the significance of results. (B) Differential binding analysis of ChIP-Seq data at three time points to monitor the activation of ER. The ER shows a strong increase in binding at 45 minutes vs. 0 minutes (C) and the majority of sites still display binding at 90 minutes.

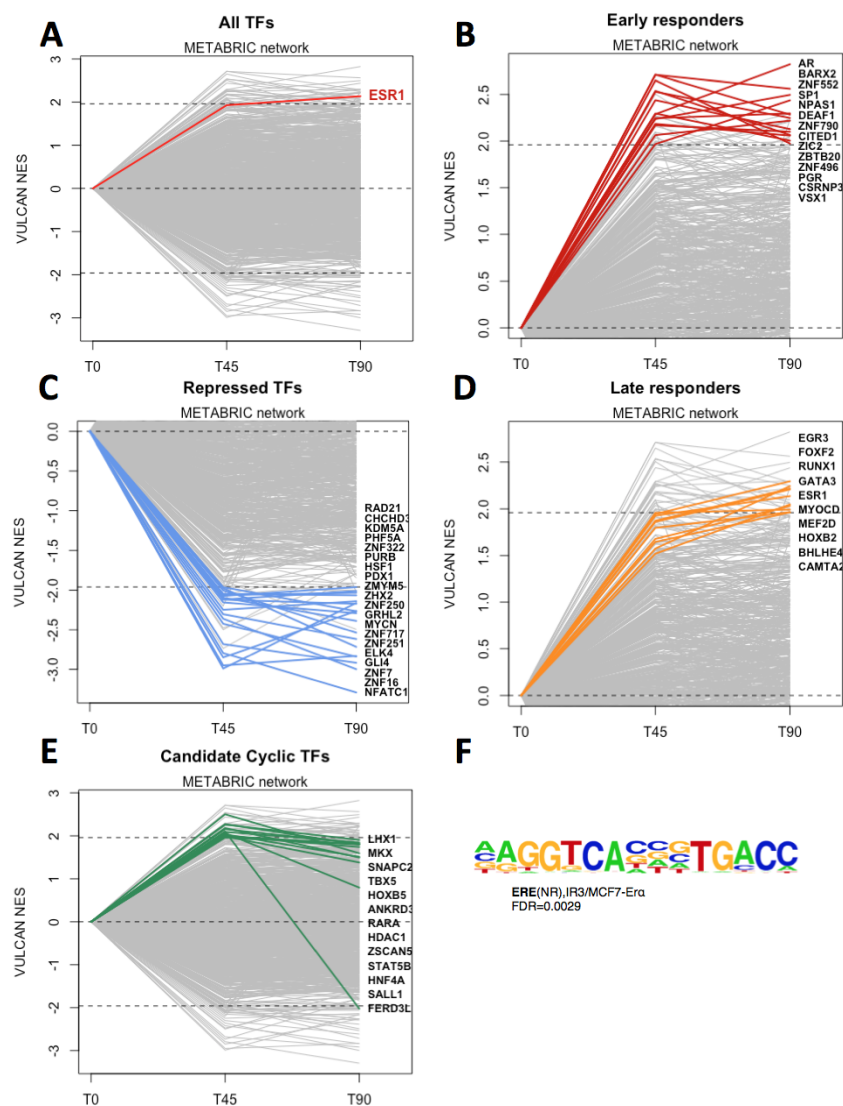
Figure 2: An overview of VULCAN



A: ChIP-Seq analysis from multiple conditions is undertaken to generate cistrome data at multiple timepoints (or conditions). B: Binding events are then compared using differential binding analysis to establish log-fold change values for individual binding events between each timepoint. C: ARACNe-AP infers all pairwise TF-target coexpression. In the example, the TCGA breast dataset is shown to infer A putative target that is correlated with ESR1. D: Minimalistic representation of the ARACNe-AP network, highlighting negative and positive regulation of targets by transcription factors. E: All the targets of a specific TF are divided in positive and negative, and tested on a differential binding signature through the msVIPER algorithm [18].



## Figure 3: ER occupancy after estradiol treatment in terms of TF network activity



**A:** Global TF network behavior as predicted by VULCAN in our ChIP-Seq dataset, highlighting the ESR1 TF at time 0 and 45/90minutes after estradiol treatment.

**B:** Global TF activity after Estradiol Treatment in MCF7 cells, inferred using the METABRIC network, highlighting TFs significantly upregulated at 45 minutes and 90 minutes

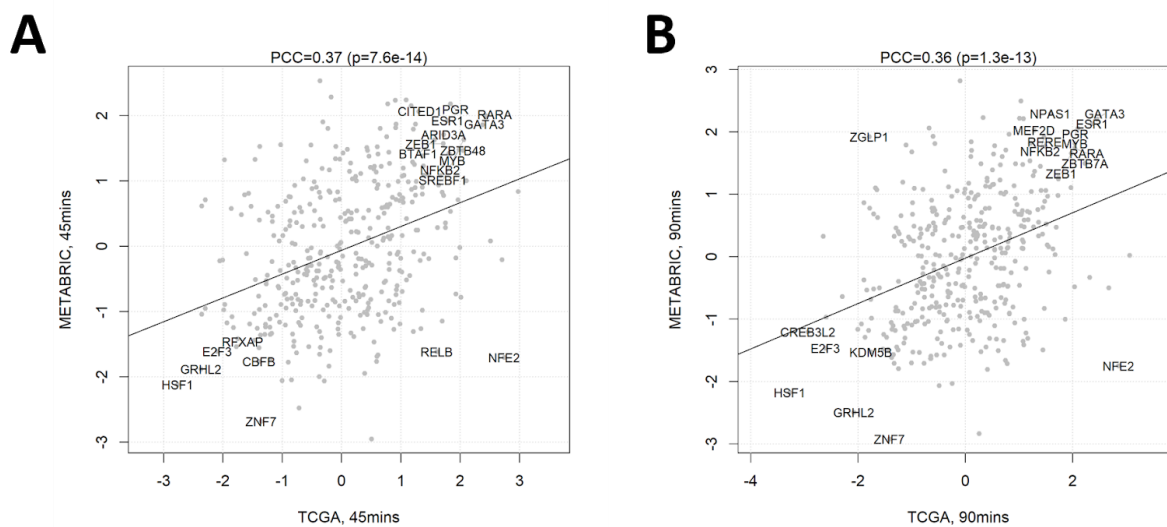
**C:** Global TF activity after Estradiol Treatment in MCF7 cells, inferred using the METABRIC network, highlighting TFs significantly downregulated at 45 minutes and 90 minutes

**D:** Global TF activity after Estradiol Treatment in MCF7 cells, inferred using the METABRIC network, highlighting TFs significantly upregulated at 45 minutes but not at 90 minutes

**E:** Global TF activity after Estradiol Treatment in MCF7 cells, inferred using the METABRIC network, highlighting TFs significantly upregulated at 90 minutes but not at 45 minutes

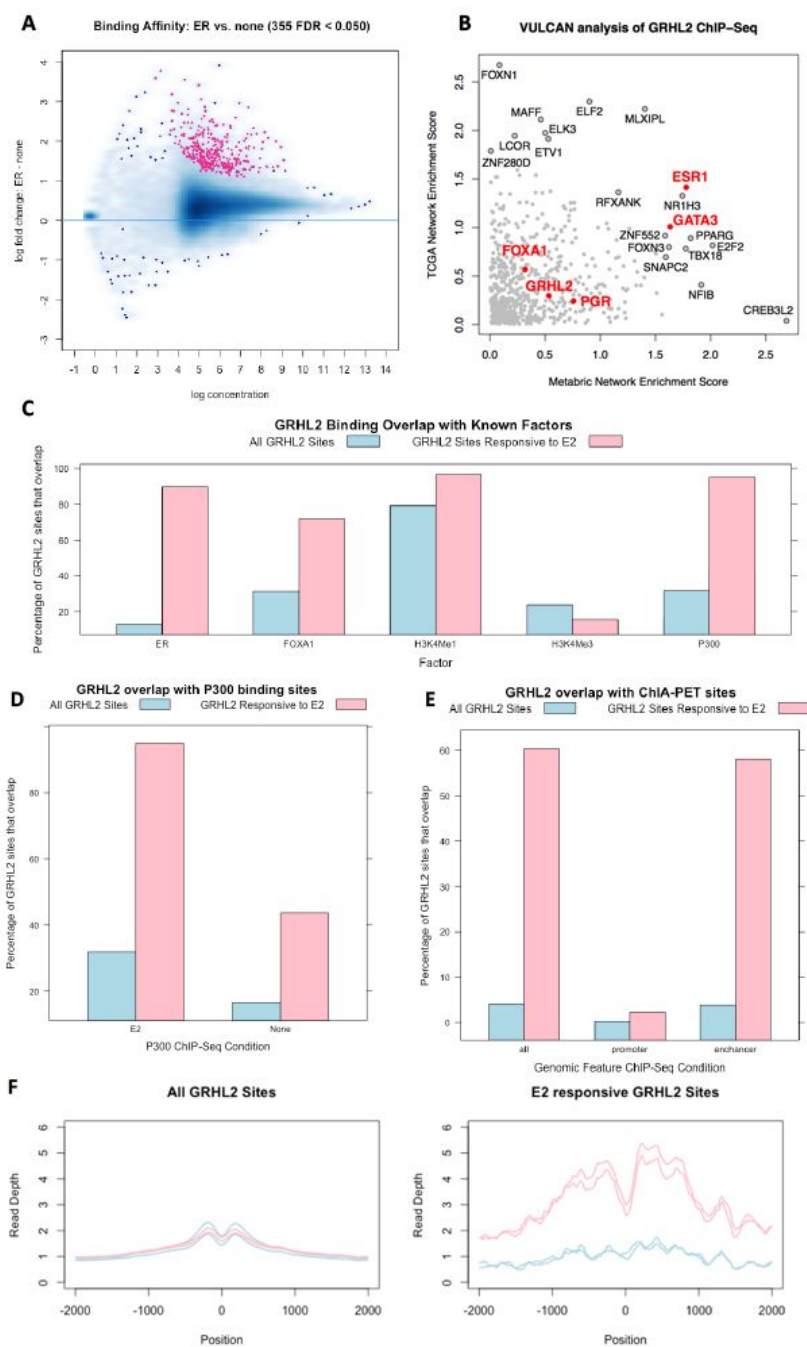
**F:** Most enriched motif in peaks upregulated at both 45 and 90 minutes after estradiol treatment, as predicted by HOMER.

Figure 4: Global TF activity after estradiol treatment using different network models.



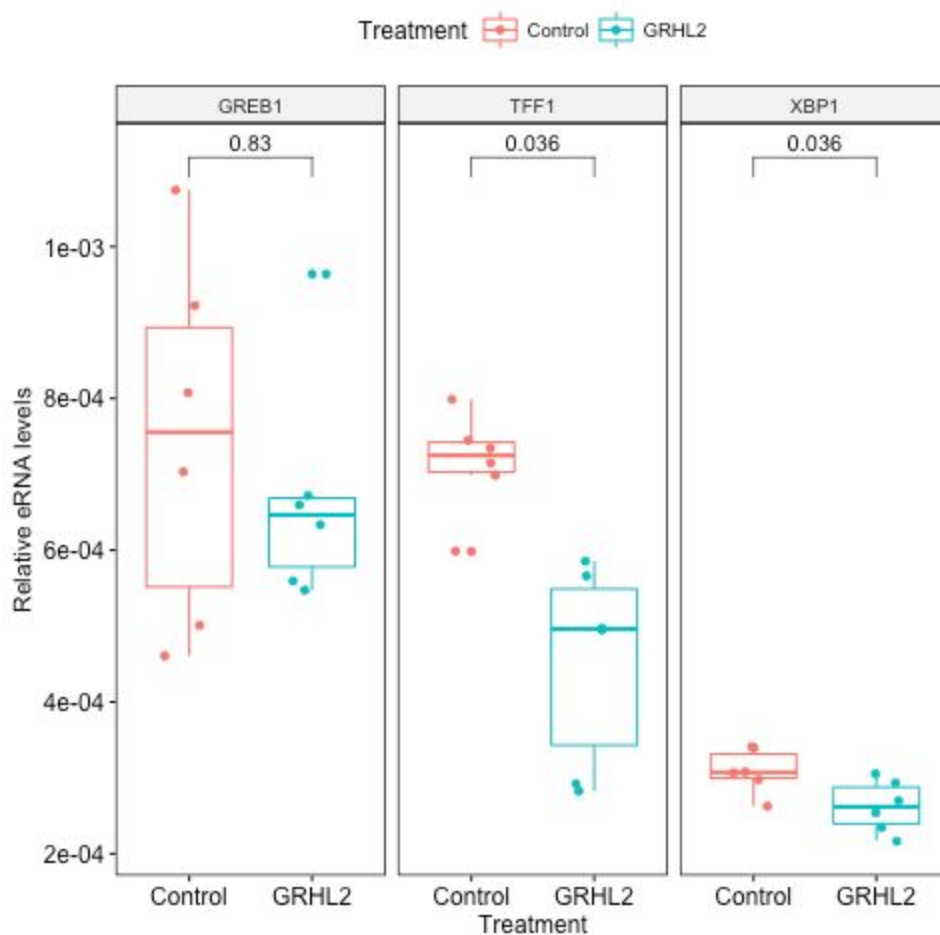
XY Scatter showing the TF activity as calculated by VULCAN for our differential ChIP-Seq analysis of ER binding at 45 minutes (A) and at 90 minutes (B) after stimulation with 100 nM E2. Comparison of results calculated using the METABRIC (y-axis) and TCGA (x-axis) networks shows consistent results know ER interactors including PGR, RARA, GATA3 and GRHL2. GRHL2 activity is notably enriched against. The regulon of ER is also consistently enriched in both networks.

## Figure 5: GRHL2 Differential ChIP-Seq between 0 and 45 minutes.



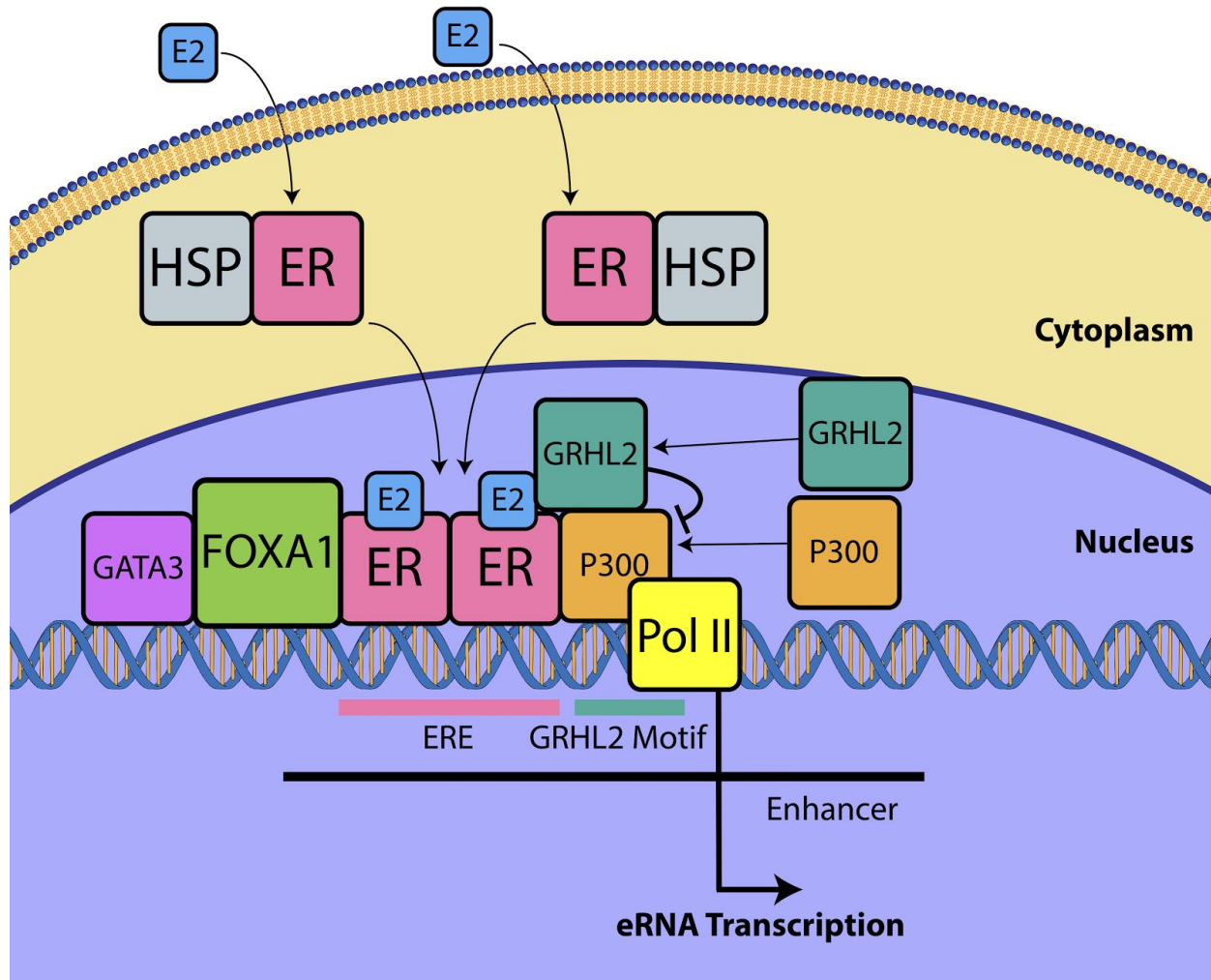
- (A)** Activation of the ER with estradiol results in a genome wide increase in GRHL2 binding.
- (B)** VULCAN Analysis of the same data show a significant enrichment for ESR1 sites in both the context of the METABRIC and TCGA networks. The regulon for FOXA1 is also not enriched. Inspection of known FOXA1/GRHL2 sites (e.g. RARα promoter) shows GRHL2 already bound.
- (C)** Overlap of GRHL2 binding with public datasets shows that E2 responsive GRHL2 sites show considerable overlap with ER, FOXA1 and P300 sites, H3K4Me1 and H3K4Me3 show little enrichment.
- (D)** Analysis of P300 binding showed a greater overlap of GRHL2 ER responsive sites in the presence of E2 than in control conditions
- (E)** Overlap with ER ChIA-PET sites showed enrichment for GRHL2 sites at ER enhancers.
- (F)** Analysis of Gro-SEQ data (GSE43836) at GRHL2 sites. Blue lines are control samples, Pink are samples after stimulation with E2. In general GRHL2 sites (left) show no change in the levels of transcription on addition of E2; however, E2 responsive GRHL2 sites (right) show a robust increase in transcription on the activation of the ER.

## Figure 6: Effect of over-expression of GRHL2 on eRNA at E2 responsive binding sites.



Overexpression of GRHL2 in MCF7 resulted in a reduction of eRNA transcribed from the GREB1, TFF1 and XBP1 enhancers. The effect was significant at TFF1 and XBP1 enhancers ( $p < 0.05$ , Wilcoxon paired-test).

**Figure 7: Overview of the role of GRHL2 in ER activation**



On activation of the ER by the ligand E2 the protein is released from a complex containing HSPs and translocates to the nucleus. The holo-ER dimer forms a core complex at Estrogen Response Elements (ERE) with FOXA1 (pioneer factor) and GATA3. ER further recruits P300 and GRHL2. GRHL2 has an inhibitory effect on P300 (a transcriptional activator interacting with TFIID, TFIIB, and RNAPII) thereby reducing the level of eRNA transcription at enhancer sites. Over-expression of GRHL2 further suppresses transcription, while knockdown of GRHL2 reverses the process.

## References:

1. Carroll JS. Mechanisms of oestrogen receptor (ER) gene regulation in breast cancer. *Eur. J. Endocrinol.* 2016;175:R41–9.
2. Lin C-Y, Vega VB, Thomsen JS, Zhang T, Kong SL, Xie M, et al. Whole-genome cartography of estrogen receptor alpha binding sites. *PLoS Genet.* 2007;3:e87.
3. Driscoll MD, Sathya G, Muyan M, Klinge CM, Hilf R, Bambara RA. Sequence requirements for estrogen receptor binding to estrogen response elements. *J. Biol. Chem.* 1998;273:29321–30.
4. Carroll JS, Meyer C a, Song J, Li W, Geistlinger TR, Eeckhoute J, et al. Genome-wide analysis of estrogen receptor binding sites. *Nat. Genet.* 2006;38:1289–97.
5. Carroll JS, Liu XS, Brodsky AS, Li W, Meyer CA, Szary AJ, et al. Chromosome-Wide Mapping of Estrogen Receptor Binding Reveals Long-Range Regulation Requiring the Forkhead Protein FoxA1. *Cell.* 2005;122.
6. Frasor J, Danes JM, Komm B, Chang KCN, Lyttle CR, Katzenellenbogen BS. Profiling of estrogen up- and down-regulated gene expression in human breast cancer cells: insights into gene networks and pathways underlying estrogenic control of proliferation and cell phenotype. *Endocrinology.* 2003;144:4562–74.
7. Zaret KS, Carroll JS. Pioneer transcription factors: establishing competence for gene expression. *Genes & Dev.* 2011;25:2227–41.
8. Horwitz KB, Jackson TA, Bain DL, Richer JK, Takimoto GS, Tung L. Nuclear receptor coactivators and corepressors. *Mol. Endocrinol.* 1996;10.
9. Glass CK, Rose DW, Rosenfeld MG. Nuclear receptor coactivators. *Curr. Opin. cell Biol.* 1997;9:222–32.
10. Green KA, Carroll JS. Oestrogen-receptor-mediated transcription and the influence of co-factors and chromatin state. *Nat. Rev. Cancer.* 2007;7:713–22.
11. Klinge CM. Estrogen receptor interaction with estrogen response elements. *Nucleic acids Res.* 2001;29:2905–19.
12. Shang Y, Hu X, DiRenzo J, Lazar MA, Brown M. Cofactor dynamics and sufficiency in estrogen receptor-regulated transcription. *Cell.* 2000;103:843–52.
13. Reid G, Hubner MR, Metivier R, Brand H, Denger S, Manu D, et al. Cyclic, proteasome-mediated turnover of unliganded and liganded ER alpha on responsive promoters is an integral feature of estrogen signaling. *Mol. Cell.* 2003;11:695–707.
14. Métivier R, Penot G, Hübner MR, Reid G, Brand H, Kos M, et al. Estrogen receptor-alpha directs ordered, cyclical, and combinatorial recruitment of cofactors on a natural target promoter. *Cell.* 2003;115:751–63.
15. Basso K, Margolin A, Stolovitzky G, Klein U, Dalla-Favera R, Califano A. Reverse engineering of regulatory networks in human B cells. *Nat Genet. Nature Publishing Group;* 2005;37:382–90.
16. Giorgi FM, Del Fabbro C, Licausi F. Comparative study of RNA-seq- and Microarray-derived coexpression networks in *Arabidopsis thaliana*. *Bioinformatics.* 2013;29.
17. Castro MAA, de Santiago I, Campbell TM, Vaughn C, Hickey TE, Ross E, et al. Regulators of genetic risk of breast cancer identified by integrative network analysis. *Nat. Genet.* 2016;48:12–21.
18. 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.

19. Lachmann A, Giorgi FM, Lopez G, Califano A. ARACNe-AP: gene network reverse engineering through adaptive partitioning inference of mutual information. *Bioinforma.* 2016;32:2233–5.
20. Curtis C, Shah SP, Chin S-F, Turashvili G, Rueda OM, Dunning MJ, et al. The genomic and transcriptomic architecture of 2,000 breast tumours reveals novel subgroups. *Nature.* 2012;486:346–52.
21. Cancer Genome Atlas Network. Comprehensive molecular portraits of human breast tumours. *Nature.* 2012;490:61–70.
22. Lefebvre C, Rajbhandari P, Alvarez MJ, Bandaru P, Lim WK, Sato M, et al. A human B-cell interactome identifies MYB and FOXM1 as master regulators of proliferation in germinal centers. *Mol. Syst. Biol.* 2010;6:377.
23. Lachmann A, Giorgi FM, Lopez G, Califano A. ARACNe-AP: gene network reverse engineering through adaptive partitioning inference of mutual information. *Bioinformatics.* 2016;32.
24. McLean CY, Bristol D, Hiller M, Clarke SL, Schaar BT, Lowe CB, et al. GREAT improves functional interpretation of cis-regulatory regions. *Nat. Biotechnol.* 2010;28:495–501.
25. Balwierc PJ, Pachkov M, Arnold P, Gruber AJ, Zavolan M, van Nimwegen E. ISMARA: automated modeling of genomic signals as a democracy of regulatory motifs. *Genome Res.* 2014;24:869–84.
26. Welch RP, Lee C, Imbriano PM, Patil S, Weymouth TE, Smith RA, et al. ChIP-Enrich: gene set enrichment testing for ChIP-seq data. *Nucleic acids Res.* 2014;42:e105.
27. Orlando DA, Chen MW, Brown VE, Solanki S, Choi YJ, Olson ER, et al. Quantitative ChIP-Seq normalization reveals global modulation of the epigenome. *Cell reports.* 2014;9:1163–70.
28. Guertin MJ, Markowitz F, Holding AN. Novel Quantitative ChIP-seq Methods Measure Absolute Fold-Change in ER Binding Upon Fulvestrant Treatment. 2017;
29. Subramanian A, Tamayo P, Mootha VK, Mukherjee S, Ebert BL, Gillette MA, et al. Gene set enrichment analysis: a knowledge-based approach for interpreting genome-wide expression profiles. *Proc. Natl. Acad. Sci. United States Am.* 2005;102:15545–50.
30. Ross-Innes CS, Stark R, Teschendorff AE, Holmes KA, Ali HR, Dunning MJ, et al. Differential oestrogen receptor binding is associated with clinical outcome in breast cancer. *Nature.* 2012;481:389–93.
31. Stein RA, Chang C-Y, Kazmin DA, Way J, Schroeder T, Wergin M, et al. Estrogen-related receptor alpha is critical for the growth of estrogen receptor-negative breast cancer. *Cancer Res.* 2008;68:8805–12.
32. Bhat-Nakshatri P, Wang G, Appaiah H, Luktuke N, Carroll JS, Geistlinger TR, et al. AKT alters genome-wide estrogen receptor alpha binding and impacts estrogen signaling in breast cancer. *Mol. Cell. Biol.* 2008;28:7487–503.
33. Gozgit JM, Pentecost BT, Marconi SA, Ricketts-Loriaux RSJ, Otis CN, Arcaro KF. PLD1 is overexpressed in an ER-negative MCF-7 cell line variant and a subset of phospho-Akt-negative breast carcinomas. *Br. J. Cancer.* 2007;97:809–17.
34. Mohammed H, D'Santos C, Serandour AA, Ali HR, Brown GD, Atkins A, et al. Endogenous purification reveals GREB1 as a key estrogen receptor regulatory factor. *Cell reports.* 2013;3:342–9.
35. Thompson A, Schäfer J, Kuhn K, Kienle S, Schwarz J, Schmidt G, et al. Tandem Mass Tags: A Novel Quantification Strategy for Comparative Analysis of Complex Protein Mixtures by



MS/MS. *Anal. Chem.* 2003;75.

36. Jozwik KM, Chernukhin I, Serandour AA, Nagarajan S, Carroll JS. FOXA1 Directs H3K4 Monomethylation at Enhancers via Recruitment of the Methyltransferase MLL3. *Cell reports.* 2016;17:2715–23.
37. Zwart W, Theodorou V, Kok M, Canisius S, Linn S, Carroll JS. Oestrogen receptor-co-factor-chromatin specificity in the transcriptional regulation of breast cancer. *EMBO J.* 2011;30:4764–76.
38. Hah N, Murakami S, Nagari A, Danko CG, Kraus WL. Enhancer transcripts mark active estrogen receptor binding sites. *Genome Res.* 2013;23:1210–23.
39. Xiang X, Deng Z, Zhuang X, Ju S, Mu J, Jiang H, et al. Grhl2 determines the epithelial phenotype of breast cancers and promotes tumor progression. *PloS one.* 2012;7:e50781.
40. Cieply B, Riley P, Pifer PM, Widmeyer J, Addison JB, Ivanov AV, et al. Suppression of the epithelial-mesenchymal transition by Grainyhead-like-2. *Cancer Res.* 2012;72:2440–53.
41. Pifer PM, Farris JC, Thomas AL, Stoilov P, Denvir J, Smith DM, et al. Grainyhead-like 2 inhibits the coactivator p300, suppressing tubulogenesis and the epithelial-mesenchymal transition. *Mol. Biol. cell.* 2016;27:2479–92.
42. Paltoglou S, Das R, Townley SL, Hickey TE, Tarulli GA, Coutinho I, et al. Novel Androgen Receptor Coregulator GRHL2 Exerts Both Oncogenic and Antimetastatic Functions in Prostate Cancer. *Cancer Res.* 2017;77:3417–30.
43. Chi D, Singhal H, Li L, Long HW, Garber JE, Brown MA. Abstract 3376: Reprogramming the estrogen signaling network is a potential mechanism for human breast tumorigenesis. *Cancer Res.* 2017;77.
44. Holmes KA, Brown GD, Carroll JS. Chromatin Immunoprecipitation-Sequencing (ChIP-seq) for Mapping of Estrogen Receptor-Chromatin Interactions in Breast Cancer. *Methods Mol. Biol.* 2016;1366:79–98.
45. Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics.* 2009;25.
46. Zhang Y, Liu T, Meyer CA, Eeckhoutte J, Johnson DS, Bernstein BE, et al. Model-based analysis of ChIP-Seq (MACS). *Genome Biol.* 2008;9:137.
47. Subramanian A, Tamayo P, Mootha VK, Mukherjee S, Ebert BL, Gillette MA, et al. Gene set enrichment analysis: a knowledge-based approach for interpreting genome-wide expression profiles: a knowledge-based approach for interpreting genome-wide expression profiles. *Proc. Natl. Acad. Sci. United States Am. National Acad Sciences;* 2005;102:15545–50.
48. Heinz S, Benner C, Spann N, Bertolino E, Lin YC, Laslo P, et al. Simple combinations of lineage-determining transcription factors prime cis-regulatory elements required for macrophage and B cell identities. *Mol. cell.* 2010;38:576–89.
49. Crooks GE, Hon G, Chandonia J-M, Brenner SE. WebLogo: a sequence logo generator. *Genome Res.* 2004;14:1188–90.
50. Anders S, Huber W. Differential expression analysis for sequence count data. *Genome Biol. BioMed Central Ltd;* 2010;11:106.
51. Hurtado A, Holmes KA, Ross-Innes CS, Schmidt D, Carroll JS. FOXA1 is a key determinant of estrogen receptor function and endocrine response. *Nat Genet. Nature Publishing Group;* 2010;43:27–33.