1  **Enhancers mapping uncovers phenotypic heterogeneity and evolution in**

2  **patients with luminal breast cancer**

3

4  Darren K. Patten[*1], Giacomo Corleone[*1], Balázs Győrffy[2,3], Edina Erdős[4], Alina

5  Saiakhova[5], Kate Goddard[6], Andrea Vingiani[7], Sami Shousha[8], Lőrinc Sándor

6  Pongor[2], Dimitri J. Hadjiminas[8], Gaia Schiavon[9], Peter Barry[10], Carlo Palmieri[11], Raul

7  C. Coombes[1], Peter Scacheri[5], Giancarlo Pruneri[12], Luca Magnani[#1].

8

9  * Equal contribution

10  [#]To whom correspondence should be addressed: l.magnani@imperial.ac.uk

11

12  [1]Department of Surgery and Cancer, The Imperial Centre for Translational and

13  Experimental Medicine, Imperial College London, Hammersmith Campus, London,

14  U.K.

15  [2]MTA TTK Lendület Cancer Biomarker Research Group, Institute of Enzymology,

16  Hungarian Academy of Sciences, 1117, Budapest, Hungary

17  [3]Semmelweis University 2nd Dept. of Pediatrics, 1094, Budapest, Hungary

18  [4]Department of Biochemistry and Molecular Biology, Genomic Medicine and

19  Bioinformatic Core Facility, University of Debrecen, Debrecen 4032, Hungary

20  [5]Department of Genetics and Genome Sciences, Case Comprehensive Cancer

21  Center, Case Western Reserve University, Cleveland, OH 44106

22  [6]Department of Breast and General Surgery, Charing Cross Hospital, Imperial

23  College Healthcare NHS Trust, London, U.K.

24  [7]Department of Pathology, European Institute of Oncology, Milan

25  [8]Centre for Pathology, Department of Medicine, Imperial, College London, Charing

26  Cross, London, UK

27  [9]Translational Science, IMED Oncology, AstraZeneca, Cambridge, UK.

1    [10]Department of Breast Surgery, The Royal Marsden NHS Foundation Trust,

2    Orchard House, Downs Road, Sutton, SM2 5PT, UK.

3    [11]Institute of Translational Medicine University of Liverpool, Clatterbridge Cancer

4    Centre, NHS Foundation Trust, and Royal Liverpool University Hospital, Liverpool,

5    Merseyside, UK

6    [12]Pathology Department, Fondazione IRCCS Istituto Nazionale Tumori and

7    University of Milan, School of Medicine.

8

9

10

1 **Abstract**

2 The degree of intrinsic and interpatient phenotypic heterogeneity and its role in

3 tumour evolution is poorly understood. Phenotypic divergence can be achieved via

4 the inheritance of alternative transcriptional programs[1,2]. Cell-type specific

5 transcription is maintained through the activation of epigenetically-defined regulatory

6 regions including promoters and enhancers[1,3,4]. In this work, we annotated the

7 epigenome of 47 primary and metastatic oestrogen-receptor (ER$\alpha$)-positive breast

8 cancer specimens from clinical samples, and developed strategies to deduce

9 phenotypic heterogeneity from the regulatory landscape, identifying key regulatory

10 elements commonly shared across patients. Highly shared regions contain a unique

11 set of regulatory information including the motif for the transcription factor YY1. *In*

12 *vitro* work shows that YY1 is essential for ER$\alpha$ transcriptional activity and defines the

13 critical subset of functional ER$\alpha$ binding sites driving tumor growth in most luminal

14 patients. YY1 also control the expression of genes that mediate resistance to

15 endocrine treatment. Finally, we show that H3K27ac levels at active enhancer

16 elements can be used as a surrogate of intra-tumor phenotypic heterogeneity, and to

17 track expansion and contraction of phenotypic subpopulations throughout breast

18 cancer progression. Tracking YY1 and SLC9A3R1 positive clones in primary and

19 metastatic lesions, we show that endocrine therapies drive the expansion of

20 phenotypic clones originally underrepresented at diagnosis. Collectively, our data

21 show that epigenetic mechanisms significantly contribute to phenotypic

22 heterogeneity and evolution in systemically treated breast cancer patients.

23

24

25

26

27

28

29

1  **Introduction**

2  Breast cancer (BC) is the most common cancer type and the second most frequent

3  cause of cancer related death in women[5]. 70% of all BC cases contain variable

4  amounts of oestrogen receptor-alpha (ERα) positive cells. ERα is central to BC

5  pathogenesis and serves as the target of endocrine therapies (ET)[6,7]. ERα-positive

6  BC is typically subdivided in two 'intrinsic' molecular subtypes (luminal A and luminal

7  B[8]) characterized by distinct prognosis, highlighting functional inter-patient

8  heterogeneity. Recent analyses demonstrate that patient-to-patient heterogeneity is

9  more pervasive (reflected by histological[9], genetic architecture[10] and transcriptional[11]

10 differences) ultimately influencing long-term response to endocrine treatment[12].

11 Indeed, 30-40% of ERα BC patients relapse during or after completion of adjuvant

12 endocrine therapies. At the time of relapse, almost all patients will have developed

13 resistance to ET, partly through treatment-specific genetic evolutionary trajectories[13].

14 Additionally, the presence of genetic intra-tumor heterogeneity has also now been

15 extensively documented in several cancer types, demonstrating the role of clonal

16 evolution in cancer [14]. Yet, recent studies have shown that driver coding-mutations

17 do not significantly change between primary and metastatic luminal breast cancer,

18 with the notable exception of *ESR1* mutations[15], suggesting that alternative

19 mechanisms might contribute to BC progression and drug-resistance. Parallel to

20 genetic evolution, phenotypic/functional changes driven by epigenetic mechanisms

21 can also contribute to breast cancer progression and ET resistance in cell lines[16,17].

22 Nevertheless, little is known about the epigenome of BC patients, its influence on

23 intra-tumour phenotypic heterogeneity and its role in breast cancer progression.

24 Epigenetic modifications consist of chemical modifications targeting the

25 DNA/RNA (e.g. DNA methylation) and the chromatin (histone modifications). Histone

26 modifications have been successfully used to map regulatory regions and to

27 annotate the non-coding DNA[1,3]. Acetylation of lysine 27 on histone 3 (H3K27ac) is

28 strongly associated with promoters and enhancers of transcriptionally active genes

29 [4,18,19]. Increasing evidence suggests that epigenetic information can actively transfer

30 gene transcription states across cell division[20-23]. Epigenetic modifications play also

31 a central role in modulating ERα binding to the DNA possibly by interacting with

32 ERα-associated pioneer factors [24,25]. Finally, *in vitro* studies have shown that

1  epigenetic reprogramming might play a central role in ERα BC cells that adapt to
2  endocrine therapy[16,17].

3      Here we show the results of a systematic investigation of the epigenetic
4  landscape of ERα positive primary and metastatic breast cancer from 47 individuals.
5  Our results represent the first large scale topographic mapping of the active
6  regulatory landscape of longitudinal ERα-positive BC. Using H3K27ac we mapped
7  active promoters and enhancers across treatment naïve primary and endocrine
8  treated metastatic patients. We used bioinformatic approaches to deconvolute the
9  complex regulatory landscape and identified inter- and intra-patient epigenetic
10 heterogeneity. We mined promoters and enhancers from clinically relevant breast
11 cancer samples for potential regulatory drivers identifying YY1 as a novel key player
12 in ERα-positive BC. Finally, we demonstrate that epigenetic mapping can efficiently
13 estimate phenotypic heterogeneity changes throughout BC progression.

14

15 **Results**

16 **Mapping enhancers and promoters of primary and metastatic ERα positive**
17 **breast cancer**

18 To build a comprehensive compendium of all the clinically relevant active regulatory
19 regions of luminal BC we profiled fifty-five ERα positive BC samples (primary n=39,
20 and metastatic n=16) with H3K27ac ChIP-seq (Supplementary Table S1). To
21 minimize the introduction of noise from non-tumor tissues we used samples with high
22 tumor burden (>70%, Supplementary Figures S1). 85% of samples yielded
23 satisfactory results (47/55, Supplementary Figure 2A and Table S2). H3K27ac-
24 enriched regions were classified into 23,976 gene-proximal (1kb upstream of
25 transcription start site (TSS), promoters) and 326,719 gene-distal (enhancers).
26 Considering the ten-fold difference in H3K27ac signal, it was not surprising to
27 observe that 80% of promoters can be captured by profiling 4 individual, while nearly
28 40 are needed to reach the same coverage for enhancers, as indicated by saturation
29 plots (Supplementary Figure 2B). These data are in agreement with enhancers being
30 the main determinants of cell-type specific transcriptional differences [4,18,26,27]. To
31 gain insights on the penetrance of each regulatory region, we developed a Sharing
32 Index (SI) by annotating all enhancers and promoters in function of the number of

patients sharing the H3K27ac signal at each specific location (Supplementary Figure 2C). In agreement with saturation analyses, we find that a large portion of enhancers are patient-specific (SI=1) while active promoters are more commonly shared between patients (Supplementary Figure 2C). Collectively, these data demonstrate that enhancers account for the majority of potential epigenetic heterogeneity in ERα-positive BC.

**Enhancer activity allows the qualitative assessment of phenotypic heterogeneity**

Genetic heterogeneity is an hallmark of most solid tumours [28]. Nonetheless, genetic intra-tumoral heterogeneity does not often directly translates to phenotypic heterogeneity. In agreement, despite extensive inter- and intratumoral clonal genetic diversity[29], the majority of ERα-positive tumors benefit from systemic ET[12]. Likewise, treatment-naïve metastatic patients generally respond well to ET, at least initially, suggesting that genetic heterogeneity on its own cannot explain treatment resistance/response. On the other hand, phenotypic hierarchies can override genetic hierarchies in brain cancers [2,30], suggesting that inheritable epigenetic program might ultimately contribute to phenotypic heterogeneity and treatment outcome.

The existence of intra-tumoral phenotypic heterogeneity in breast cancer patients has been known to pathologists for decades, at least for a small number of biomarkers. For example, immunohistochemistry (IHC) assessment of the proportion of ERα-positive within luminal cancer patients varies on a *continuum* from less than 1% to nearly 100%[31]. Unfortunately, IHC is low-throughput method, and typically only a few proteins can be studied before the sample is consumed. In contrast, assessing phenotypical heterogeneity from patient bulk transcriptional data is unfeasible as transcription is ultimately an analogue signal in which each individual cell can contribute a stochastic amount of RNA, making data deconvolution impractical (Fig 1A). For instance, cells with focal gene amplification have higher bulk gene expression but individual cells can contribute radically different amounts as shown by single-molecule single-cell RNA FISH[13]. On the other hand, recent evidence show that the signal captured by chromatin assays such as ATAC-seq appears to be directly proportional to the cells contributing to it[32]. Similarly, ChIP-seq signal can be

1  thought of as digital information with each single nucleosome being ON (K27ac) or

2  OFF at any given time (Fig. 1A). Notably, even within genetically clonal cell lines, the

3  H3K27ac signal varies considerably between different regulatory regions. Regulatory

4  regions labelled as super enhancers, for example, have 10-100-times more

5  H3K27ac signal than typical enhancers[18]. What accounts for the variation in signal is

6  not known, but one possibility is that heterogeneity within the cell population (either

7  clonal or sub-clonal) contribute to signal intensity. For example, super-enhancers

8  might represent regulatory regions active across most or all cells within a population

9  at any given time (clonal, H-peaks), while "typical" enhancers with lower H3K27ac

10 signal may represent sub-clones (M and L peaks, Fig. 1A). This concept is similar to

11 measuring variant allele frequencies (VAF) to infer genetic heterogeneity.

12 Phenotypical heterogeneity might be the consequence of heterogeneous cell

13 populations (i.e. tumor, stroma and immune infiltrate) or actual cancer-specific

14 epigenetic subclones. As our ChIP-seq data are derived from high tumor burden

15 samples, we hypothesized that H3K27ac signal could allow for a qualitative

16 assessment of phenotypic heterogeneity. We further theorised that direct correlation

17 might exist between clonal prevalence (intra-tumour) and population prevalence

18 (inter-patients) (Fig. 1A).

19      We tested the initial assumption by performing spike-in experiments in which

20 known numbers of cells with well-characterized regulatory region activity (and similar

21 genetic background)[16] were admixed in incremental proportions prior to H3K27ac

22 ChIP-qPCR. The data shows that H3K27ac enrichment is proportional to the number

23 of cells with the active enhancer (Fig. 1B). These findings are corroborated by an

24 independent analysis using a different antibody (ERα) (Supplementary Figure S3A).

25 As the signal between different patients is not directly comparable, we normalized

26 the data using a ranking approach, assigning to each H3K27ac signal a Rank Index

27 (RI, 1 to 100, strongest to weakest) (Fig. 1C). Signal from low RI (H peaks) might be

28 associated with clonal regulatory regions active in almost all cells. Conversely, high

29 RI (M-L peaks) mark more heterogeneous/sub-clonal enhancer activity. By

30 investigating the relationship between RI and SI we found extremely high correlation

31 between these two parameters (Fig. 1D), suggesting that clonal regulatory regions

32 are more common between patients while sub-clonal regulatory elements are more

33 patient-specific. We defined clonal low-RI/high-SI loci as regulatory drivers (RD,

34 SI>21) and high-RI/low-SI that might originate from sub-clonal populations as

regulatory noise (RN, SI<21). Of note, a small but discrete proportion of promoters/enhancers escape this general trend having extremely low RI despite being patient-specific or higher RI while being shared (dot-boxes, Fig. 1D).

**Enhancers are associated with BC risk-SNP and control gene transcription**

We next investigated the extent to which regulatory regions identified in our cohort associate to BC. Previous analyses from ERα-BC cell lines, have shown that genetic predisposition might occur through SNPs that modulate transcription factors binding at enhancers (FOXA1 and ERα[33]). We then tested the relationship between DNA risk variants specifically associated with BC through GWAS[33,34] and regulatory regions captured in patients. Strikingly, almost the totality of validated BC risk variants is contained within our H3K27ac database. Currently, this dataset represents the most enriched annotation for GWAS variants in breast cancer (Fig 1E). This overlap is highly significant specifically for enhancers but not for other annotations (Fig. 1F). Notably, this association is not replicated using colorectal cancer risk variants suggesting that these enhancers might play a specific role in BC development (Fig. 1F).

Next, we assessed the relationship between estimated enhancers clonality and transcriptional output. Transcriptional data obtained using microarray and RNA-seq estimate the average expression level within a population. The average expression is function of the number of cells engaged in active transcription and the number of RNA molecule within each cell[35]. Interestingly, several lines of evidence suggest that RNA transcription is stochastic [13,36], thus implying that the total number of cells with active transcription significantly contribute to changes in average RNA levels in bulk populations. As our analysis allows for qualitatively prediction of clonality in enhancer activity, it allows to test if clonal enhancers active in the majority of cells correlate with higher RNA levels. To do so we linked enhancers to their potential target genes using CTCF insulated boundaries[37]. We then analysed three independent BC expression datasets (METABRIC [10], TCGA [38] and Affymetrix [39]) in function of RI/SI indexes. Our analyses show a predictable increase for mRNA levels with parallel increases in the associated SI, further suggesting that RDs drive RNA expression in a progressively increasing number of cells (Supplementary Figure 3B). These results were more modest when analysing the transcriptome from normal

1    breast tissue (Supplementary Figure 3B, small insets) suggesting that our analysis

2    has identified a subset of regulatory regions strongly associated with malignant

3    outgrowth. These data indicate that transcripts identified as dis-regulated in BC

4    might reflect changes in the size of phenotypic subpopulations. This could be in part

5    driven by a selection process, as normal breast transcriptional data reflect the

6    heterogeneous composition of the tissue (adipocyte, myoepithelial and epithelial

7    cells), while BC is normally dominated by epithelial morphology. Collectively, our

8    data show that enhancer activity strongly tracks transcriptional changes in breast

9    cancer patients.

10

## Imputed transcription factors landscape of ERα breast cancer patients

12    Enhancers stores regulatory information in the form of transcription factors (TFs)

13    binding motifs[40]. The vast majority of TFs require accessible chromatin in order to

14    bind their cognate DNA sequences [41]. We reasoned that a systematic investigation

15    of the predicted TFs landscape in function of enhancer activity could reveal potential

16    transcriptional BC drivers. To narrow down to the accessible DNA within active

17    enhancers and promoters we integrated the DNaseI signal from 129 cell lines with

18    the inferred nucleosome pattern obtained from the H3K27ac data (Fig. 2A). Initial

19    analyses collapsing all imputed DHS in relationship with their enhancers and

20    promoter location identified correctly well-known BC-TFs according to their

21    promoter–enhancer bias (Supplementary Figure 4). We then stratified the complete

22    set of enhancers and promoter regions based on the associated SI and repeated TF

23    motif analysis focusing within each SI-defined bin followed by unsupervised

24    clustering. This analysis generated two major clades (Fig. 2B), indicating the

25    presence of different classes of regulatory regions. Strikingly, we find that RD and

26    RN enhancers and promoters cluster specifically into the two major clades,

27    suggesting that putative clonal and sub-clonal enhancers contain distinct regulatory

28    information (Fig. 2B). Functional TF binding is associated with TF leaving a footprint

29    within chromatin accessible regions [40,42]. Interestingly, clonal enhancers in ERα-

30    positive MCF7 breast cancer cells are significantly enriched in TF footprints[16], while

31    sub-clonal enhancers are significantly deprived of footprints suggesting that TFs

32    might bind clonal enhancers with longer residence time [42] (Fig. 2C).

1       We then focused on estrogen-response elements (ERE) as they constitute the
2   canonical DNA sequence to which ERα binds. Unexpectedly, ERE motifs are
3   enriched only in RN enhancers, suggesting that a significant amount of ERα binding
4   occurs in sub-clonal/less functional enhancers (Fig. 2B). To gain further insights on
5   ERα dynamics we turned to a recently published ERα dataset obtained from patient
6   material (n=15) [43]. Generally, the proportion of ERα binding sites overlapping
7   enhancers increase with the SI (9% vs. 70%, SI-1 vs. SI-39, Fig. 2C). This was not
8   observed for promoters (15% vs. 6%, SI-1 vs. SI-39, Fig. 2C), and is consistent with
9   previous studies demonstrating a significant bias for ERα binding at active enhancer
10  elements[44,45]. These data imply that shared enhancers have a strong propensity for
11  ERα binding despite being generally under-enriched in EREs (Fig. 2B). More
12  importantly, the bulk of ERα binding were captured only once in fifteen patients (ERα
13  SI=1), with less than 0.003% of ERα being shared across 75% of the patients (484
14  core ERα)[43] (Fig. 2D). Together, these data support biochemical evidence that
15  suggest that only a small fraction of ERα binding events with longer-residency time is
16  functional[42]. We therefore concluded that the largest portion of ERα binding identified
17  in patients occur at patient specific, sub-clonal enhancers and might be the
18  consequence of the transient ERα-DNA interactions occurring while the receptor
19  scans the genome[42]. The discrepancy between the small number functional ERα
20  core binding and the observation of ERE-poor RD enhancers led us to hypothesize
21  that other TFs might collaborate with ERα to increase its transcriptional efficiency at
22  clonal enhancer. Most TF motifs enriched in RD enhancers are also largely observed
23  in RN regions, with the notable exception of the motif for the transcription factor YY1
24  (Fig. 2B). Interestingly, TF analysis of the footprints within MCF7 clonal enhancer
25  (Fig. 2C, RI<20) similarly identifies YY1 as the top hit (qVal=0.001). YY1 has been
26  recently implied in *de novo* formation of enhancer promoter looping during neural
27  development [46] and MYC-like ability to potentiate gene expression[47] indicating a
28  potential role in modulating the enhancer landscape in ERα-positive BC.

29

30  **YY1 enhancer activity mark a dominant phenotypic clone in BC**

31  YY1 is a ubiquitously expressed TF (Supplementary Figure 5A-B) that can act as an
32  activator or repressor by binding DNA, RNA and chromatin modifiers[48,49]. YY1

1   function in breast cancer is poorly understood, as it has been linked to different

2   outcomes depending on breast cancer cells subtypes. In luminal BC, YY1 appears to

3   be positively correlated with AP2 to promote HER2 activity [50], while in triple negative

4   BC models appears to be a tumor suppressor by controlling BRCA1 expression [51].

5   Interestingly, YY1 drosophila homolog PhoRC is involved in epigenetic memory by

6   recruiting of Polycomb repressor complex to sequence specific regions[52], but YY1's

7   role in mammals is not entirely understood [46,47,53]. Our TF analysis shows that YY1

8   might actually operate as a global reader of active clonal enhancers common to the

9   majority of ERα-positive patients. These data therefore predict that most luminal

10  breast cancers should contain a dominant YY1-positive clone. To assess the size of

11  YY1 phenotypic clone in our patient's dataset we identified the *bona fide* enhancers

12  looping at YY1 promoter using 3D chromatin maps [54] (Supplementary Fig. 6A). We

13  found three potential enhancers with high SI within a CTCF-insulated region with

14  YY1 promoter (SI A=41, B=33 and C=26, Supplementary Fig. 6A). Interestingly

15  Enhancer A also directly interacts with Enhancer B-C, suggesting a multi-enhancers

16  interaction with YY1 promoter. Enhancer A consistently ranks among the most clonal

17  enhancer in nearly all patients, suggesting that YY1 is transcribed in almost all cells

18  (Fig. 3A). By comparison, in most normal tissues profiled by H3K27ac within the

19  Epigenome Roadmap consortium[45], YY1 Enhancer A activity is more variable while

20  remaining relatively dominant, implying that some tissues may harbour YY1-

21  subclonal subpopulations (Fig 3B). Consistent with these predictions,

22  immunocytochemistry (IHC) meta-analysis (Fig 3B) showed a decreasing number of

23  YY1 positive cells in correspondence to increasing RI scores (Insets, Fig 3B and

24  Supplementary Figure 6B). Collectively, these data suggest that enhancer ranking

25  can capture qualitative changes in intra-tumoral heterogeneity, and that YY1-

26  enhancer activity marks a dominant phenotypic clone in ERα-positive BC.

27  Next, we looked at the significance of YY1 mRNA expression in a pan-cancer

28  analysis and found that tumor tissue generally have significantly higher expression

29  level for YY1 as compared to normal tissues (Supplementary Figure 7A). This does

30  not appear to reflect simply the proliferation status as we found no correlation

31  between YY1 and Ki67 expression in 2509 Breast Cancer patients [55]. This

32  observation was replicated in an independent large BC dataset as well (Fig. 3C and

33  Supplementary Figure 7B). Of note, YY1 is not subject to recurrent genomic

1  aberrations (data not show). These data suggest that cancer lesions might contain a

2  larger fraction of YY1-expressing cells as compared to more heterogeneous tissues

3  (Fig. 3B). Meta-analysis of the METABRIC[10] datasets shows that patients with higher

4  YY1 mRNA level at diagnosis have significantly worse outcome (Fig. 3C). We

5  replicated this finding using TCGA RNA-seq data with significant stratification

6  occurring in Luminal A patients, a subtype typically associated with good prognosis

7  (Fig. 3C). To test if YY1 increased mRNA expression could be driven by an

8  expansion of YY1-positive cells from a more heterogeneous population, we stained

9  normal breast section with IHC. Our data show that lobules and ducts contain distinct

10 YY1 positive sub-clonal populations within the luminal and basal compartments (Fig

11 3D). On the other hand, nearby tumour tissue is overwhelmingly YY1 positive,

12 demonstrating the existence of a clonal YY1 population. The data demonstrate that

13 the expansion of a YY1 phenotypic clone drive the changes in bulk RNA levels

14 between normal and tumor samples and reinforce the notion that YY1 might play a

15 central role in ERα-positive BC.

16

**YY1 is a global enhancers modulator marking functional ERα binding**

18      The TF analysis in our epigenomic patient dataset revealed YY1 uniquely

19 associated with clonal regulatory regions and potentially collaborating with ERα at

20 critical enhancers. To gain mechanistic insight on the role of YY1 we performed YY1

21 ChIP-Seq in quiescent (estrogen-deprived) and estrogen-stimulated luminal BC

22 MCF7 cells. Cells were stimulated with estrogen for 45 minutes, upon which

23 maximum ERα-binding to chromatin occurs[45]. ChIP-seq biological replicates show

24 very high correlation ($R^2$=0.98), thus we kept consensus loci for further analyses. In

25 quiescent cells, YY1 occupies a very small set of enhancers and promoters near

26 housekeeping genes[56] (Fig. 4A). Strikingly, estrogen stimulation induce a 23-fold

27 expansion of the YY1 binding repertoire, mostly at enhancer regions (Fig. 4A). Newly

28 occupied loci are associated with ERα-BC signatures and epigenetic editors (Fig.

29 4A). Interestingly, only ~10% of all binding is characterized by a high affinity YY1

30 motif suggesting that induced YY1 could also bind directly to modified nucleosomes

31 through its chromatin remodelling partner INO80[57]. Orthogonal analyses show that

32 induced-YY1 binding involves almost all MCF7 active regulatory regions and is

strongly associated with H3K27ac marks (Fig. 4B). Nonetheless, it is unlikely that YY1 binding directly promote or requires H3K27ac as we did not find any difference between quiescent (estrogen-deprived) and estrogen-stimulated H3K27ac epigenomes (data not shown). Conversely, YY1 binding is absent from silenced genes (Supplementary Figure 7C), demonstrating that YY1 does not associate with PRC2 mediated repression in BC cancer cells.

Our *in vivo* analysis suggest that clonal YY1-bound enhancers are generally not enriched for EREs or ERα, with the exception of the atypical core-ERα (Fig 2B). In agreement, our *in vitro* data show only marginal overlap between YY1 and ERα or its pioneer factor FOXA1 (Fig. 4B-D) indicating that generally YY1 recruitment is independent from ERα. On the other hand, YY1, ERα and FOXA1 co-localize at increased frequencies at core-ERα loci in MCF7 cells (Fig. 4C). Similar observations were made by comparing YY1 overlap with *in vivo* derived ERα binding (60% overlap with core ERα vs. 18% overlap with patient-unique ERα). In addition, we find that genes defining the luminal subtype in TCGA patients are significantly enriched for ERα core binding with YY1 but not patient-unique ERα (Fig. 4D). Overall, these data further suggest that YY1 might stabilize ERα binding[42] at a small subset of transcriptionally productive enhancers (core-ERα) captured in most tumor cells and most patients[43]. To test if YY1 can contribute to ERα driven transcription, we measured luciferase activity from a promoter driven by an array of estrogen response elements (EREs) in MCF7 cells in presence or absence of YY1(Fig. 4E) and show absolute dependencies on YY1. Furthermore, YY1 depletion also abrogates cell proliferation in response to estrogen stimulation in MCF7 (Fig. 4F) suggesting that YY1 is a direct driver of the clonal proliferation observed in the BCa (Fig. 3D-E). These observations were replicated in other independent luminal BC cell models (ZR75 and T47D, Supplementary Figure 8A-B). Finally, we show that YY1 depletion leads to significant downregulation of core-ERα target genes in luminal BC cell line models (Supplementary Figure 8C). Collectively these data identify YY1 as a novel essential transcription factor significantly contributing to ERα regulatory network transcriptional activity.

**YY1 contributes to drug-resistance in luminal BC**

1    YY1 motif is highly enriched in clonal enhancers identified in primary and metastatic

2    luminal patients (Fig 2B). All metastatic patients included in this study relapsed

3    following adjuvant endocrine therapies suggesting that YY1 might also play role in

4    this setting. In agreement, primary and metastatic samples show clonal YY1

5    enhancer activity, indicating that YY1 positive cells are not effectively cleared by the

6    therapy (Fig. 3A). Therefore, we investigated the role of YY1 in LTED cells, an

7    MCF7-derivative that develop estrogen-independent growth partly through

8    constitutive activation of ERα signalling[16]. YY1 depletion leads to complete

9    abrogation of LTED growth demonstrating that YY1 is still required at this stage (Fig.

10   4G). Interestingly, LTED cells have an expanded repertoire of ERα binding

11   compared to MCF7, fuelled by endogenous ligands [13,16]. The set of enhancers

12   engaged by ERα and YY1 in LTED cells is radically different compared to MCF7,

13   with the majority of ERα-YY1 being specific to each cell type (Fig 4I, LTED only:

14   3598/5037). ERα-YY1 bound enhancers in LTED strongly associates with the

15   transcription of genes involved with acquired endocrine therapy, suggesting that

16   during epigenetic reprogramming[16], YY1 might stabilize ERα to LTED specific

17   enhancers (Fig. 4J). To further examine the relationship between YY1 and endocrine

18   resistance we analyzed a set of estrogen responsive genes whose transcription

19   cannot be antagonized by Tamoxifen treatment in MCF7 cells[58]. These genes were

20   not enriched for patient-private ERα, but we saw an ever-increasing association with

21   ERα-YY1 bound enhancers, especially with core ERα-YY1 (Fig. 4K). For example,

22   ERα-YY1 is found near CXXC5 and SLC9A3R1, ranked respectively first and second

23   as the most strongly estrogen-induced gene that cannot be antagonized by

24   Tamoxifen [58]. Collectively, these data strongly support the role of YY1 in ERα BC

25   growth and progression.

26

27   **YY1-ERα promote SCL9A3R1 expression despite endocrine treatment**

28   SLC9A3R1 (NHERF1/EBP50)[59] SLC9A3R1 encodes a Na/H exchanger regulatory

29   cofactor. SLC9A3R1 null mice have disrupted protein-kinase-A-dependent cAMP-

30   mediated phosphorylation[60]. In agreement with SLC9A3R1 potential role in

31   endocrine resistance, meta-analysis of patient-derived data using all available genes

32   (n=22,277) reveals that SLC9A3R1 expression is amongst the top 1% of genes with

the strongest prognostic association with relapse in a cohort of 724 ET-treated ERα-positive patients [39] (Fig. 5A). High expression of SLC9A3R1 also significantly correlates with poor survival in additional independent ERα-BC datasets (Supplementary Figure 9A). In addition to Tamoxifen treatment, SLC9A3R1 remains transcriptionally active in most endocrine therapy resistant BC cell lines that retain ERα expression (Supplementary Figure 9B) but genetic or pharmacological (Fulvestrant) suppression of ERα is sufficient to block SLC9A3R1 transcription (Supplementary Figure 9B-C). Specifically, SLC9A3R1 expression is not antagonized by estrogen deprivation[61] (LTED models, Supplementary Figure 9B-D), nor Raloxifene *in vitro* (Supplementary Figure 9E) or neo-adjuvant AI treatment in clinically treated patients *in vivo* (Fig. 5B). Overall, these data demonstrate that SLC9A3R1 is a direct ERα target whose expression cannot be antagonized by first-line endocrine therapies.

Bulk RNA-seq data from a panel of cancer cell lines demonstrate that ERα-positive BC cells have the highest levels of SLC9A3R1 mRNA (Supplementary Figure 10A). More importantly, TCGA RNA-seq analysis shows that SLC9A3R1 expression is higher specifically in ERα BC patients compared to normal tissue or other subtypes (Supplementary Figure 10B). Chromatin analyses of MCF7 and LTED cells identify ER-bound enhancers at 3 independent loci within the insulated SLC9A3R1 locus (E1-E3), a RD region directly looping to the YY1-bound SLC9A3R1 promoter within a CTCF insulated perimeter (Supplementary Figure 10C). Strikingly, E1 and E2 contain core ERα bindings in addition to YY1, while E3 contains a patient unique ERα binding with no YY1 (Supplementary Figure 10C). *In vivo* transcriptional analysis demonstrates that SLC9A3R1 is the only gene near the E1-E2 enhancers that shows a significant increase in bulk-RNA level when comparing normal breast tissue with ERα–positive BC (Supplementary Figure 10D). Interestingly, enhancer-activity appears to be immune to endocrine therapy (MCF7 vs. MCF7 tamoxifen resistant and LTED, Supplementary Figure 10C). Collectively, these data strongly support the notion that SLC9A3R1 expression is driven by a breast cancer specific enhancer within the expanding ERα-YY1 clone during tumor initiation. Nonetheless, SLC9A3R1 expression is dependent on YY1 (Supplementary Figure 11A), demonstrating that both ERα and YY1 are essential for full enhancer activity. Silencing SLC9A3R1 is sufficient to abrogate oestrogen-induced growth in ERα-

1    positive cells (Fig 5C). Intriguingly, SLC9A3R1 is not essential for a second in ERα-

2    positive model (T47D) but appears to be a critical gene for both AI-resistant cells

3    models (Fig. 5C and Supplementary Figure S11B). Collectively, these data

4    demonstrate that ERα-YY1 regulate SLC9A3R1 via enhancer binding and identify

5    SLC9A3R1 as a novel player involved in ET resistance.

6

7    **Mapping phenotypic heterogeneity using YY1 and SLC9AR1 enhancer activity**

8    Both SLC9A3R1 and YY1 enhancers are commonly activated in our patient's dataset

9    (SI=34 and SI=41 respectively). Yet, YY1 enhancer identifies YY1-positive cells as a

10   dominant clone in almost all patients (RI≤20, Fig 3A). Conversely, SLC9A3R1

11   enhancer activity indicates that SLC9A3R1 marks a potentially dynamic sub-clonal

12   population in most primary patients (RI≥20, Fig 5D). Our *in vitro* data suggest that

13   SLC9A3R1 transcription cannot be antagonized by endocrine therapies while

14   SLC9A3R1 is important for resistant BC cell lines. This predicts that the SLC9A3R1-

15   positive population should increase under adjuvant treatment. Interestingly, the only

16   evidence of a clonal SLC9A3R1 population was found in samples from three

17   metastatic, endocrine-resistant patients (Fig. 5D).

18        Bulk transcriptional data show that average SLC9A3R1 expression is

19   significantly higher in ERα positive BC cells but do not inform about potential

20   subpopulations (Supplementary Figure 10A-B). Conversely, meta-analysis of

21   SLC9A3R1 enhancer activity (RI) within the ENCODE H3K27ac datasets indicates

22   that MCF7 cells contain a clonal SLC9A3R1 population while all other cell lines

23   appear to have decreasing sub-clonal populations (Supplementary Figure 11C). Of

24   note, the size of the sub-clonal population correlates with total RNA content for the

25   cells contained in both assays, suggesting that the decreasing bulk RNA signal is

26   driven by a progressively smaller subpopulation (Supplementary Figure 11C). Similar

27   analysis of YY1 enhancer indicate that cancer cell lines are vastly clonal for YY1

28   expression (Supplementary Figure 11D) in agreement with clinical samples. Notably,

29   both YY1 and SLC9A3R1 RIs in mammary epithelial cells suggest the presence of

30   sub-clonal populations. These observations fit well with experimental data from IHC

31   profiles from normal breast (Fig. 3D and Supplementary Figure 12B). Meta-analyses

32   of H3K27ac from the Epigenome Roadmap database predict that most tissues

1   potentially have only sub-clonal populations, as determined by SLC9A3R1 enhancer

2   activity (Fig. 5E). In agreement, the size of the sub-clonal population tracks the

3   mRNA signature of each tissue with SLC9A3R1 potentially clonal tissues

4   accumulating the highest amount of RNA-seq tags (Small and Large Intestine,

5   Supplementary Figure 12). Analogously to YY1, meta-analysis of IHC data identifies

6   decreasing SLC9A3R1-positive with increasing RI scores (Fig. 5E and

7   Supplementary Figure 12B). Finally, to further validate that RI index can estimate

8   phenotypic clones, we retrospectively collected available FFPE biopsies for the BC

9   patients profiled with H3K27ac ChIP-seq (n=19). We then performed IHC using YY1

10  (Fig. 5F) and SLC9A3R1 (Fig. 5G) antibodies and compared the predicted enhancer

11  activity (RI) with the actual size of YY1 and SLC9A3R1-positive clones within each

12  patient. With the exception of one metastatic sample (M3), YY1 staining robustly

13  correlate with RI, confirming large clonal YY1 positive populations in all examined

14  tissues (Fig. 5F). In parallel, SLC9A3R1 enhancer activity correctly estimated the

15  size of the sub-clonal subpopulations in individual patients (Fig 5D). Further meta-

16  analyses on Protein Atlas data support these findings, by identifying YY1 clonal

17  populations and SLC9A3R1 sub-clonal populations in most ERα BC samples

18  (Supplementary Figure 13A-C). Overall, these data strongly support the notion that

19  enhancer activity can be used to qualitatively deconvolute heterogeneous

20  populations into phenotypical subclones.

21

22  **Phenotypic evolution during BC progression is shaped by endocrine treatment**

23  Tumor evolution studies have primarily focused on treatment naïve patients, taking

24  advantage of multi-regional sampling to retrospectively monitor changes in

25  clonality[14,62]. Clonal tracking is dependent in part on passenger mutations, and the

26  effect of therapy has been rarely accounted for[13,63]. More importantly, clonality has

27  been traced uniquely using genetic variants, with the intrinsic limitation of correlating

28  genetic changes to phenotypic ones. For example, sub-clones defined by passenger

29  mutations might be phenotypically equivalent, while a recent study using barcoded

30  glioblastoma cells shows that phenotypic clones might evolve independently from

31  mutational signatures [2]. In addition, the few studies that looked at driver mutations in

32  coding regions of primary and metastatic BC disease found relatively similar

17

1  mutational landscapes[15], suggesting that mapping phenotypic clones though BC

2  progression might reveal new targets. Our ability to acquire qualitative estimates of

3  phenotypic clones using enhancer ranking provides for a potential approach for

4  tracking changes in tumor heterogeneity with the additional advantage of predicting

5  for potentially functional changes. We interrogated our patient's dataset focusing on

6  events occurring between treatment-naïve primaries and treatment-resistant

7  metastatic BC (Fig. 6A). We hypothesized that phenotypic clonal evolution might be

8  driven by a coordinated activation/selection of groups of enhancers during BC

9  progression and this could be influenced by treatment. Our previous results suggest

10  that YY1+ cells remain clonal during progression (Fig 3A). Conversely, we show that

11  SLC9A3R1 expression is not antagonized by endocrine treatment suggesting that

12  SLC9A3R1-positive clones could expand during progression. We then calculated

13  changes in RI ($\Delta$RI) for all enhancers captured in at least three patients (SI>3,

14  n=88935) between primary and metastatic samples (Fig. 6B). SLC9A3R1 ranks

15  amongst the enhancers with the strongest increase in predicted clonality going from

16  primary to metastatic samples (Fig. 6B-C, 3.86$\sigma$ from median $\Delta$RI). Conversely, YY1

17  enhancer activity remains relatively unchanged (Fig. 6B-C). These data support our

18  initial hypothesis and suggest that SLC9A3R1-positive clones might expand in

19  response to treatment. To substantiate these data, we mapped the size of YY1 and

20  SLC9A3R1 phenotypic clones using IHC in an independent series of 20 matched

21  longitudinal biopsies. All surgical biopsies were obtained from treatment naïve

22  patients, while all the metastatic biopsies were taken at first relapse after endocrine

23  treatment[13]. We found YY1+ cells clonal both in primary and metastatic biopsies

24  (Fig. 6D). Conversely, SLC9A3R1+ subclones significantly expand during metastatic

25  progression to become completely clonal (100% staining) in 13/20 patients.

26  Interestingly, the only metastatic case in which we have observed a contraction of

27  the SLC9A3R1+ clone also showed a concomitant loss of ERα and PR positivity,

28  confirming our *in vitro* analysis and demonstrating that SLC9A3R1 remains an ERα

29  dependent-target despite being ET insensitive *in vivo* (Fig. 6D, red line and

30  Supplementary Figure 9B-E). Overall, these data demonstrate that changes in

31  enhancer ranking can estimate functional evolution during breast cancer

32  progression.

To gain more insight on functional evolution, we systematically annotated all regulatory regions based on bias in detection between primary and metastatic patients (Fig 6E). As expected, the bulk of enhancers and promoters do not show bias toward primary and metastatic BC patients (common enhancers, CE). However, we could successfully identify two distinct sets of regulatory regions that are preferentially associated with primary (primary enhancers, PE) or metastatic (metastatic enhancers, ME) patients. Remarkably, while CE do not show stage-specific changes in RI, PE underlie larger sub-clonal populations in primary cancers (statistically higher RIs in primary compared to metastatic, Fig. 6E). Likewise, ME have lower RI in metastatic samples suggesting that the number of cells carrying these enhancers have increased during progression (Fig. 6E). We next explored the potential causes and functional consequences driving these coordinated epigenetic changes. We thus identified the potential transcriptional targets of our enhancers taking in account CTCF boundaries[37]. Strikingly, we find that PE-associated gene-transcription is associated with significantly better outcome while ME-associated gene-transcription in primary samples is associated with poor prognosis (Fig 6D). These data imply that primary samples containing larger subpopulations of phenotypic clones with metastatic features relapse earlier.

We then mined PE and ME regulatory regions to identify the associated biological features[56]. PE appear to promote abnormal proliferation and vascularization, two key events in early tumorigenesis. Remarkably, metastatic samples switch to functional clones characterized by genes associated with BC progression (FOXA1[43]) or endocrine therapy resistance[64,65] (Fig. 6E). Altogether, these data suggest that endocrine therapies play a central role in shaping phenotypic clonal evolution. Additional in-depth studies are needed to dissect the temporal events triggered during phenotypic clonal evolution. Phenotypic subclones could evolve by early coordinated activation and decommissioning of epigenetically defined regulatory regions (*acquired*), selection of the fittest pre-existent epigenomic landscape (*de novo*) or a combination of both.

**Discussion**

1  Our work describes the first systematic epigenetic profiling of primary and metastatic
2  luminal breast cancer and reveals several critical principles underlying phenotypic-
3  functional heterogeneity and its role in breast cancer progression. By mapping
4  H3K27ac in untreated and treated patient samples we have also identified YY1 and
5  SLC9A3R1, two new key players contributing to BC. While genomic profiling of
6  breast cancer patients has revealed extensive clonal heterogeneity and retrospective
7  tumour evolution[28,66], the vast majority of the mutational burden can be considered
8  composed of passenger mutations[29,67] making difficult to extrapolate actual
9  phenotypes. Most RNA-based analysis, which may better reflect the phenotypic state
10  of cancer cells, is generally obtained from bulk tissue and cannot inform on the
11  existence of distinct subpopulations. Finally, molecular pathology can inform on the
12  relative amount of protein abundance at the single-cell level but is laborious and not
13  suitable for testing multiple targets simultaneously. In this work, we used epigenomic
14  analyses to extrapolate phenotypic heterogeneity in solid tumour samples. Our
15  analysis reveals that histone-based ChIP-seq signals, similarly to ATAC-seq[32],
16  generally correlates with the number of cells in a population carrying the specific
17  epigenetic information. Our predictions using YY1 and SLC9A3R1 enhancer fit
18  extremely well with experimental data derived from normal tissues or BC patients.
19  The findings that clonal regulatory regions dominating the landscape of individual
20  tumor samples are shared across many patients, parallel recent genomic evidences
21  showing that truncal (high allele frequency) mutations are also the most common
22  mutations within cancer cohorts.

23        The results described here have several practical implications for BC. First, by
24  comparing samples from drug-resistant metastatic patients with drug-naïve primary
25  samples, we uncovered a set of enhancers marking phenotypic clones that
26  significantly expand during breast cancer progression. Notably, these enhancers are
27  strongly associated with genes specifically transcribed in cells that acquire endocrine
28  therapy resistance (Fig. 6H). Conversely, enhancers progressively lost during tumour
29  progression are linked to processes that often occur early in tumorigenesis. A set of
30  enhancers expanding in metastatic samples point at progressive activation of
31  FOXA1 and its network. It was recently reported that FOXA1 levels are increased in
32  metastatic samples[43,68]. Our data then predict that, similarly to SLC9A3R1, FOXA1
33  positivity increases as a consequence of the expansion of a phenotypic clone

marked by an active FOXA1 enhancer and not via increased transcription of the FOXA1 gene within single cells. It is tempting to speculate that this paradigm might be valid for other genes. If correct it might signify that during cancer evolution, the proportion of cells activating transcription is more important than the absolute changes in transcription at the single cell levels. Interestingly, a set of enhancers deactivated during progression involve IL-2 signalling (Fig. 6H). Reduction in IL-2 signalling was identified as a potential marker of relapse[69]. Whether the IL-2 signal source is the BC cells[70] themselves or it is due to a small contamination of immune cells, needs to be defined. Equally, it will be important to measure real-time activation/selection of enhancers in appropriate systems to ultimately establish if phenotypic cancer evolution can be driven by Lamarckian events.

Finally, our analysis has identified two novel drivers of luminal BC. Firstly, we identified YY1 as a key TF associated with clonal enhancers and promoters in BC patients. Our data strongly support the idea that YY1 acts as a global co-activator in cancer cells associating with the entire active epigenetic landscape. Several lines of evidence indicate that YY1 might interact directly with modified nucleosomes, possibly through its partner INO80[57]. YY1 widespread association with clonal enhancer suggests it might play a role in epigenetic memory. Intriguingly, a positive screen for factors that improve induced pluripotent cells formation (iPS), identified YY1 as the top hit, further supporting its potential role as enhancer gatekeeper [71]. More specifically to ERα BC, we hypothesize that YY1 plays a critical role to stabilize ERα binding at the transcriptionally productive core- ERα enhancers. Single-molecule imaging shows that estrogen activated ERα increases its residency time on the chromatin[42] and recent evidence has shown that eRNA can trap YY1 on the chromatin [49]. More importantly, enhancer co-occupancy for YY1 and ERα occurs almost exclusively at highly shared-highly functional core-ERα bound loci. Altogether, these data raise the intriguing hypothesis that YY1 might contribute to increased ERα residency at clonal enhancers. This could explain why some ERα are captured in most patients, as longer residency time would increase chances of being captured by ChIP-Seq[43]. Longer residency might also explain the increased transcriptional activity (Fig. 4D) and increased TF footprints (Fig. 2C) of these enhancers.

YY1-ERα jointly control SLC9A3R1 enhancer activity, an event that cannot be antagonized by conventional first-line endocrine therapies (Tamoxifen or AI) and that drives SLC9A3R1 clonal expansion during breast cancer progression. Of note, primary patients with high SLC9A3R1 expression might be viewed as containing larger population of SLC9A3R1+ cells, thus resembling drug-resistant BC. Intriguingly, SLC9A3R1 is amongst the strongest single prognostic genes for relapse-free survival when considering endocrine treated patients (Fig. 5A). An attractive possibility is that YY1 stabilizes ERα sufficiently at the SLC9A3R1 enhancer maintaining epigenetic memory in the presence of external antagonists. Future studies are required to investigate the exact mechanisms through which SLC9A3R1 contribute to BC and efficient strategies to antagonize its transcription, possibly using CDK4/6 inhibitors[72] to destabilize YY1. We recently demonstrated that individual endocrine therapies can drive parallel genetic evolution *in vivo* [13] and epigenetic reprogramming *in vitro*[16]. Our data now strongly support the notion that therapeutic interventions also play an essential role driving specific epigenetic evolution during BC progression in the clinic. Metastatic re-biopsy at the time of relapse, which is becoming commonplace in clinical practice should then examine epigenetic changes in addition to newly acquired genomic ones, especially when no new genetic drivers to guide further treatment are apparent[15]

**On Line Materials and Methods**

**Tumour tissue processing**

Breast cancer sample for ChIP-seq were collected by Imperial Tissue Bank (project ethic approval R15021) and from Breast Cancer Now Tissue Bank (BCNTB-TR000053-MTA & TR000040). Breast cancer fresh frozen tissue samples each undergo aseptic macroscopic adipose tissue dissection. The dissected tumour tissue is sectioned into 2mm x 2mm fragments in a petridish placed over dry ice. Tumour fragments are then fixed using 1% formaldehyde solution for 10 minutes. Cold glycine (1M) is added to the formaldehyde-fixed tissue for 10 minutes. The tumour fragments are then pulverised using pestle and mortar and homogenised using liquid nitrogen.

**Chromatin immunoprecipitation (ChIP)**

The ChIP protocol was conducted as described by Schmidt et al.[73] with few modifications. In summary, following fixation, the tumour tissue undergoes chromatin extraction and sonication using the Bioruptor Pico sonication device (Diagenode; B01060001) using 20 cycles (30s on and 30s off) at maximum intensity. Purified chromatin was then separated for 1. Immunoprecipitation using 4ug of H3k27ac antibodies (Abcam; ab4729) per ChIP experiment or using 4ug of YY1 antibodies (Santa Cruz; sc-281 X). ChIP-seq experiment for YY1 were performed in biological duplicates. 2. Non-immunoprecipitated chromatin, used as Input control and 3. Assessment of sonication efficiencies using a 1% agarose gel. Before construction of ChIP-seq libraries (NEB Ultra II kit, see supplementary methods), enrichment of the immunoprecipitated sample was ascertained using positive and negative controls for ChIP-qPCR. Library preparation was performed using 10 – 50 ng of immunoprecipitated and Input samples.

**ChIP-qPCR**

1   Briefly, reactions were carried out in 10 ul volume containing 5 ul of Sybergreen mix

2   (ABI; 4472918), 0.5 ul of primer (5 uM final concentration), 2.5 ul of genomic DNA

3   and 2 ul of DNASE/RNASE –free water. A three-step cycle programme and a

4   melting analysis were applied. The cycling steps were as follows: 10s at 95 oC, 30s

5   at 60 oC and 30s at 72 oC, repeated 40 times.

6

7   **Ranking and Sharing Index**

8   See Supplementary Computational Methods.

9

10  **VSE**

11  See Supplementary Computational Methods.

12

13  **DHS imputations and TF motif analyses**

14  See Supplementary Computational Methods.

15

16  **Imputed DHS with vivo ERα binding Overlap**

17  ERα binding from in breast cancer patients were obtained from [43]. ERα sharing index

18  was calculated as before (see Supplementary Computational Methods). Overlap with

19  imputed DHS was calculated using BedTools calculating the overlap (at least one

20  base       pair)      via      Cistrome      Pipeline      Analysis      Suite

21  (http://cistrome.org/Cistrome/Cistrome_Project.html). The percentage of overlap

22  were calculated using binned DHS as variable first dataset and all the concatenated

23  in vivo ERα as second dataset.

24

25  **Footprint analysis**

26  See Supplementary Computational Methods.

27

1 **Encode and Epigenomic Roadmap Ranking**

2 See Supplementary Computational Methods.

3 **Immunocitochemistry**

4 Hematoxylin and eosin staining of clinical samples was performed to calculate tumor

5 burden prior to ChIP-seq. Briefly, 4-µm-thick sections were obtained from formalin-

6 fixed and paraffin-embedded specimens. After de-waxing in xylene and graded

7 ethanol, sections were incubated in 3% H2O2 solution for 25 minutes to block

8 endogenous peroxidase activities and then subjected to microwaving in EDTA buffer

9 for antigen retrieval. For YY1 (Protein Atlas HPA001119, Atlas Antibodies

10 Cat#HPA001119, RRID:AB_1858930) the flowing conditions were used: tissue

11 sections were incubated with the primary monoclonal. overnight at 4°C, and

12 chromogen development was performed using the Envision system (DAKO

13 Corporation, Glostrup, Denmark). A minimum of 500 tumor cells were scored with

14 the percentage of tumor cell nuclei in each category recorded. For SLC9A3R1

15 (HPA9672 and HPA27247, Atlas Antibodies Cat#HPA009672, RRID:AB_1857215

16 and Atlas Antibodies Cat#HPA027247, RRID:AB_10601162 respectively) the

17 following conditions were used. HPA9672 was diluted 1:400 and HPA27247 was

18 diluted 1:1500. Staining was automatized with a Ventana Benchmark Ultra using

19 epitope retrieval ER2 for 20 minutes. ER and PgR immunoreactivity was assessed

20 by the FDA-approved ER/PR PharmDX kit (Dako). The prevalence of ER/PgR

21 positive invasive cancer cells, independent of their staining intensity, was

22 quantitatively annotated in the original reports. In accordance with ASCO/CAP

23 guidelines, tumours with ≥1%of immunoreactivity were considered positive

24

25 **Cell culture**

26 MCF7 was cultured using Dulbecco's modified Eagle's medium (DMEM) containing

27 10% fetal calf serum (FCS) and 100 U penicillin/0.1 mg ml$^{-1}$ streptomycin, 2mM L-

28 glutamine plus $10^{-8}$ 17-β-estradiol (SIGMA E8875). MCF7 long term oestrogen

29 deprived (MCF7-LTED) cells were grown in phenol-free DMEM with 10% charcoal-

30 stripped FCS (DCFCS) and 100 U penicillin/0.1 mg ml$^{-1}$ streptomycin and 2mM L-

31 glutamine. T47D and T47D-LTED cells were passaged using DMEM containing 10%

1   FCS and 100 U penicillin/0.1 mg ml$^{-1}$ streptomycin, 2mM L-glutamine and phenol-

2   free DMEM with 10% DCFCS and 100 U penicillin/0.1 mg ml$^{-1}$ streptomycin and

3   2mM L-glutamine, respectively. ZR75-1 cells were grown in DMEM containing 10%

4   FCS and 100 U penicillin/0.1 mg ml$^{-1}$ streptomycin, 2mM L-glutamine.

5

6   **sIRNA**

7   Small interfering RNA (siRNA) against SLC9A3R1 (Gene ID; 9368: Ambion; s17919,

8   s17920), YY1 (Gene ID; 7528: Ambion; s14958, s14959, s14960) and *Silencer*

9   negative control (Ambion; AM4611). 1.5 x 10$^5$ cells were seeded per well using a 6-

10  well plate. MCF7 cells were seeded in phenol-free DMEM with 10% DCFCS and 100

11  U penicillin/0.1 mg ml$^{-1}$ streptomycin and 2mM L-glutamine. Following 24 hours, the

12  cells were then transfected with siRNA using Lipofectamine 3000 (Invitrogen;

13  L3000015). T47D and ZR75-1 cells were seeded in DMEM containing 10% FCS and

14  100 U penicillin/0.1 mg ml$^{-1}$ streptomycin, 2mM L-glutamine. Following 24 hours, the

15  cells were then transfected with siRNA using Lipofectamine 3000 (Invitrogen;

16  L3000015). Cells were harvested for analysis following at least 48 hours of

17  transfection.

18

19  **Cell lysis and Western blot**

20  Cells were washed twice in ice-cold PBS and lysed in RIPA (Sigma-Aldrich; R02780)

21  buffer supplemented with protease (Roche 11697498001) and phosphastase

22  (Sigma-Aldrich 93482) inhibitors for 30 minutes with intermittent vortexing. Samples

23  were centrifuged at 4$^o$C at maximum speed for 30 minutes after which, the

24  supernatant is transferred to a clean Eppendorf. Protein concentrations for each

25  sample was ascertained using the bicinchoninic acid (BCA) assay (ThermoFisher

26  Scientific; 23227). Equal amounts of lysates were loaded into BOLT 4-12% Bis-Tris

27  Plus Gel (Invitrogen; NW04120BOX). Proteins were transferred to a Biotrace

28  nitrocellulose membrane (VWR; PN66485) and incubated with primary antibodies

29  overnight. Proteins were then visualised using goat anti-mouse (ThermoFisher

30  Scientific; 31446) and anti-rabbit (ThermoFisher Scientific; 31462) HRP conjugated

31  secondary antibodies. Amersham ECL start Western Blotting Detection reagent (GE

1    Healthcare Life Sciences; RPN3243) was used for chemiluminescent imaging using

2    the Fusion solo (Vilber; Germany) imager.

3

**Transcriptional profiling**

5    Following 48 hours of transfection, MCF7 cells were either treated with $10^{-8}$ 17-β-

6    estradiol (SIGMA E8875) or control treatment for 6 hours prior to RNA extraction.

7    T47D and ZR75-1 cells lines were harvested for RNA following 48 hours of

8    transfection. No treatments were added.

9

**RNA extraction and real-time PCR**

11    Total RNA was extracted using RNeasy Mini Kit (Qiagen; 74106), and the cDNA was

12    reverse transcribed from 1ug of RNA using iScript cDNA synthesis kit (Bio-Rad;

13    #1708891). Real time-qPCR (RT-qPCR) reactions were carried out in 10 uL volume

14    containing 5 uL of sybergreen mix (ABI; 4472918), 0.5 ul of primer (2.5 uM final

15    concentration), 2.5 ul of genomic DNA and 2 ul of DNASE/RNASE–free water. A

16    three-step cycle programme and a melting analysis were applied. The cycling steps

17    were as follows: 10s at 95 $^{\circ}$C, 30s at 60 $^{\circ}$C and 30s at 72 $^{\circ}$C, repeated 40 times[19].

18

**Luciferase reporter assay**

20    MCF7 cells were seeded in a 24-well plates at $5 \times 10^{4}$ cells per well in phenol-free

21    DMEM with 10% DCFCS and 100 U penicillin/0.1 mg ml$^{-1}$ streptomycin and 2mM L-

22    glutamine. After 24 hours of incubation, transfection of plasmid DNA was performed

23    using Lipofectamine 3000 (Invitrogen; L3000015). Cells were transfected with 100ng

24    of ERE_Luciferase reporter, 10ng of the renilla luciferase control plasmid (pRL-

25    CMV), 10ng of pSG5_ER-α, 15 nm of siRNA and 280ng of Bluescribe DNA (BSM)

26    per well; totalling 400ng of DNA/well. After 12 hours of transfection the media was

27    replaced with fresh phenol-free DMEM with 10% DCFCS and 100 U penicillin/0.1 mg

28    ml$^{-1}$ streptomycin and 2mM L-glutamine. Treatment with $10^{-8}$ 17-β-estradiol (SIGMA

29    E8875) or control treatment was administered and the cells incubated for 24 hours.

30    Cell lysates are then obtained using Passive lysis 5X buffer (Promega; E1941). The

1    firefly and renilla luciferase activity was determined using DualGlo luciferase assay

2    kit (Promega; E2920) according to the manufacturer protocol. The renilla luciferase

3    activity measurement was utilised as control for transfection efficiency and therefore

4    the ERE_Luciferase activity was normalised to the reading obtained for the renilla

5    luciferase activity.

6

**SRB assay**

8    Briefly, the sulphorhodamine B (SRB) assay was used to monitor the effects of

9    silencing either SLC9A3R1 or YY1, using siRNAs, on cell proliferation monolayer

10    cultures. Cells were seeded in flat-bottomed 96-well plates (Costar; CLS3585) at a

11    density of $2 \times 10^3$. Cells were allowed to attach overnight after which, the first plate

12    (Day 0) is assayed after the cells have become adherent. Prospective plates are

13    assayed sequentially after 3 days, 5 days and 7 days. The cells are fixed by adding

14    200uL of cold 40% (weight/volume) of trichloroacetic acid (TCA) to each well for at

15    least 60 minutes. The plates were washed five times with distilled water and then

16    100 uL/well of SRB (0.4% wt/vol SRB in 1% wt/vol acetic acid) reagent is added to

17    each well and the plates are allowed to incubate for 30 minutes. The plates were

18    then washed five times in 1% (wt/vol) acetic acid and allowed to dry overnight. SRB

19    solubilisation was performed by adding 100 uL/well of 10 mM Tris HCl to the plates

20    and allowed to shake for 30 minutes. Optical density was then measured using the

21    Sunrise microplate reader (Tecan; Sunrise) at 492 nm. Cell proliferation is then

22    calculated over the 7-day period using Day 0 as a baseline measurement.

23

**Enrichment scores.**

25    Overlap for ER$\alpha$ (*in vivo*) vs enhancers and promoters were calculated by betoold

26    intersect were the percentage overlap is calculated over the total number of

27    regulatory regions within each bin against the concatenate ER$\alpha$ binding set (all ER$\alpha$

28    in all patients). For YY1, FOXA1 and ER$\alpha$ in MCF7, intersections were calculated

29    using Cistrome. YY1 BEDFILEs were the consensus narrow peaks of two biological

30    experiment, FOXA1 ChIP-seq data and ER$\alpha$ were obtained in house[16]. The core

31    ER$\alpha$ was BEDFILE was obtained by converting the published dataset from [43] to

1  HG19. The private ER$\alpha$ BEDFILE was obtained by iterative process to identify ER$\alpha$

2  binding unique to single patients prior to concatenation into a single file. Overlap

3  represent the fraction of the original datasets (first dataset) overlapping with core

4  ER$\alpha$ (second dataset). The TCGA luminal signature was obtained from [38]. Each

5  gene was extended for 20Kb upstream keeping in consideration the direction of

6  transcription. A null gene list was generated by subtracting the TCGA luminal

7  signature from a genome-wide gene list. Genes from the null list were extended in a

8  similar way and enrichment was calculated by comparing the fraction of TCGA gene

9  list with nearby binding vs. the null list. A list of estrogen target genes that do not

10  respond to Tamoxifen was obtained from [58]. Each gene was extended for 20Kb

11  upstream keeping in consideration the direction of transcription. A null gene list was

12  generated by subtracting the signature from a genome-wide gene list. Genes from

13  the null list were extended in a similar way and enrichment was calculated by

14  comparing the fraction of TAM resistant estrogen dependent gene list with nearby

15  binding vs. the null list.

16

17  **RI-IHC correlation**

18  FFPE sections for the patients used in the ChIP-seq section were retrieved from

19  Imperial Tissue bank. Sections were stained with YY1 or SLC9A3R1 antibodies.

20  Stained sections were divided in 20 sectors. 5 sectors with high tumor burden were

21  scored for the number of IHC+ cells and results averaged. The number of IHC+ cells

22  and the matched RI was analyzed using linear regression using Prism 5 (GraphPad

23  software Inc.).

24

25  **$\Delta$RI**

26  See Supplementary Computational Methods.

27

28  **YY1 and SLC9A3R1 Pan cancer expression analysis**

29  YY1 and SLC9A3R1 expression profile for matched Normal vs. Cancer samples was

30  obtained using TIMER diff.exp option (https://cistrome.shinyapps.io/timer/). YY1

1  transcriptional analyses of breast cancer subtypes was performed in the Metabric

2  Dataset (Curtis Breast) using probe ILMN_1770892 or TCGA dataset using

3  Oncomine (https://www.oncomine.org/resource/login.html).

4

## SLC9A3R1 Meta-analyses

6  SLC9A3R1 expression profile in drug resistant cell lines was performed by analysis

7  of RNA-seq data from [16]. SLC9A3R1 expression profile in MCF7 cells transfected

8  with siRNA against ERα was performed by analysis of microarray data from

9  GSE27473. SLC9A3R1 expression profile in additional LTED models was performed

10  by analysis of microarray data from E-GEOD-19639. All statistical analyses were

11  performed using Prism 5 (GraphPad software Inc.). Kaplan-Meier analysis using

12  SLC9A3R1 expression were performed by re-analysis of 23 independent microarray

13  datasets (KMPLOT), TCGA RNA-seq data or the combined Metabric Dataset.

14  SLC9A3R1 transcriptional profile in breast cancer cell lines was obtained from the

15  HPA RNA-seq dataset (http://www.proteinatlas.org/about/download). SLC9A3R1

16  transcriptional profile from tissues was obtained from the HPA, GTEx and FANTOM5

17  RNA-seq datasets (http://www.proteinatlas.org/about/download).

18

## Author Contribution

20  L.M. conceived the study. D.P., E.E performed the experiments. L.M., G.C., B.G.,

21  A.S., L.S., and P.S., developed and performed bioinformatic analyses. K.G.,

22  organised tissue collection. D.H., G.S., P.B., C.P., R.C.C., recruited patients and

23  supplied tissues. S.S., performed pathology assessment of ChIP-seq processed

24  samples. G.P., provided matched material. A.V. and G.P., performed IHC staining

25  and scoring. All authors read and approved the manuscript.

26

## Acknowledgments

28  We want to acknowledge and thanks all patients and their families for the support

29  and for donating the research samples. We thank Breast Cancer Now Tissue Bank

30  (project TR0121), Imperial Tissue Bank and the LEGACY study for contributing

1   tissues. The authors gratefully acknowledge infrastructure support from the Cancer

2   Research UK Imperial Centre, the Imperial Experimental Cancer Medicine Centre

3   and the National Institute for Health Research Imperial Biomedical Research Centre.

4   L.M was supported by a CRUK fellowship (P64250) and Imperial Junior Fellowship

5   (G53019). D.P was supported by a Wellcome Trust PhD studentship

6   (103034/Z/13/Z). G.C was supported by a Marie Skłodowska Curie Training Grant

7   (642691, EpiPredict). We acknowledge Lorna Watson, Iros Barozzi, Ylenia Perone

8   and Jason Carrol for their constructive comments on the manuscript.

9

10  **Data Accession**

11  H3K27ac data for all patients samples have been deposited at the ENA

12  (http://www.ebi.ac.uk/ena) under project number XXXXXX.

13

14  **References**

15  1.  Roadmap Epigenomics Consortium *et al.* Integrative analysis of 111 reference
16      human epigenomes. *Nature* **518,** 317–330 (2015).
17  2.  Lan, X. *et al.* Fate mapping of human glioblastoma reveals an invariant stem
18      cell hierarchy. *Nature* **60,** 5954 (2017).
19  3.  Consortium, T. E. P. *et al.* An integrated encyclopedia of DNA elements in the
20      human genome. *Nature* **488,** 57–74 (2012).
21  4.  Ernst, J. *et al.* Mapping and analysis of chromatin state dynamics in nine
22      human cell types. *Nature* **473,** 43–49 (2011).
23  5.  Ferlay, J. *et al.* Cancer incidence and mortality worldwide: sources, methods
24      and major patterns in GLOBOCAN 2012. *Int. J. Cancer* **136,** E359–86 (2015).
25  6.  Carroll, J. S. Steroids, nuclear receptors and breast cancer. Preface. *Mol. Cell.*
26      *Endocrinol.* **382,** 623 (2014).
27  7.  Ali, S., Buluwela, L. & Coombes, R. C. Antiestrogens and Their Therapeutic
28      Applications in Breast Cancer and Other Diseases. *Annu. Rev. Med.* **62,** 217–
29      232 (2011).
30  8.  Perou, C. M. *et al.* Molecular portraits of human breast tumours : Article :
31      Nature. *Nature* **406,** 747–752 (2000).
32  9.  Genestie, C. *et al.* Comparison of the prognostic value of Scarff-Bloom-
33      Richardson and Nottingham histological grades in a series of 825 cases of
34      breast cancer: major importance of the mitotic count as a component of both
35      grading systems. *Anticancer Res.* **18,** 571–576 (1998).
36  10. Curtis, C. *et al.* The genomic and transcriptomic architecture of 2,000 breast
37      tumours reveals novel subgroups. *Nature* **486,** 346–352 (2012).
38  11. Koboldt, D. C. *et al.* Comprehensive molecular portraits of human breast
39      tumours. *Nature* **490,** 61–70 (2012).
40  12. Early Breast Cancer Trialists' Collaborative Group (EBCTCG). Aromatase

31

1    inhibitors versus tamoxifen in early breast cancer: patient-level meta-analysis
2    of the randomised trials. *Lancet* (2015). doi:10.1016/S0140-6736(15)61074-1
3    13.    Magnani, L. *et al.* Acquired CYP19A1 amplification is an early specific
4    mechanism of aromatase inhibitor resistance in ERα metastatic breast cancer.
5    *Nature Genetics* **49,** 444–450 (2017).
6    14.    McGranahan, N. & Swanton, C. Biological and therapeutic impact of intratumor
7    heterogeneity in cancer evolution. *Cancer Cell* **27,** 15–26 (2015).
8    15.    Yates, L. R. *et al.* Genomic Evolution of Breast Cancer Metastasis and
9    Relapse. *Cancer Cell* **32,** 169–184.e7 (2017).
10   16.    Nguyen, V. T. M. *et al.* Differential epigenetic reprogramming in response to
11   specific endocrine therapies promotes cholesterol biosynthesis and cellular
12   invasion. *Nature Communications* **6,** 10044 (2015).
13   17.    Magnani, L. *et al.* Genome-wide reprogramming of the chromatin landscape
14   underlies endocrine therapy resistance in breast cancer. *Proceedings of the
15   National Academy of Sciences* **110,** E1490–E1499 (2013).
16   18.    Whyte, W. A. *et al.* Master transcription factors and mediator establish super-
17   enhancers at key cell identity genes. *Cell* **153,** 307–319 (2013).
18   19.    Heintzman, N. D. *et al.* Distinct and predictive chromatin signatures of
19   transcriptional promoters and enhancers in the human genome. *Nature
20   Genetics* **39,** 311–318 (2007).
21   20.    Falahi, F. *et al.* Towards Sustained Silencing of HER2/neu in Cancer By
22   Epigenetic Editing. *Molecular Cancer Research* **11,** 1029–1039 (2013).
23   21.    Laprell, F., Finkl, K. & Müller, J. Propagation of Polycomb-repressed chromatin
24   requires sequence-specific recruitment to DNA. *Science* eaai8266 (2017).
25   doi:10.1126/science.aai8266
26   22.    Wang, X. & Moazed, D. DNA sequence-dependent epigenetic inheritance of
27   gene silencing and histone H3K9 methylation. *Science* eaaj2114 (2017).
28   doi:10.1126/science.aaj2114
29   23.    Coleman, R. T. & Struhl, G. Causal role for inheritance of H3K27me3 in
30   maintaining the OFF state of a Drosophila HOX gene. *Science* eaai8236
31   (2017). doi:10.1126/science.aai8236
32   24.    Magnani, L., Eeckhoute, J. & Lupien, M. Pioneer factors: directing
33   transcriptional regulators within the chromatin environment. *Trends Genet.* **27,**
34   465–474 (2011).
35   25.    Jozwik, K. M. & Carroll, J. S. Pioneer factors in hormone-dependent cancers.
36   1–5 (2012). doi:10.1038/nrc3263
37   26.    Hnisz, D. *et al.* Convergence of developmental and oncogenic signaling
38   pathways at transcriptional super-enhancers. *Molecular Cell* **58,** 362–370
39   (2015).
40   27.    Heintzman, N. D. *et al.* Histone modifications at human enhancers reflect
41   global cell-type-specific gene expression. *Nature* **459,** 108–112 (2009).
42   28.    Yates, L. R. *et al.* Subclonal diversification of primary breast cancer revealed
43   by multiregion sequencing. *Nature Medicine* **21,** 751–759 (2015).
44   29.    Williams, M. J., Werner, B., Barnes, C. P., Graham, T. A. & Sottoriva, A.
45   Identification of neutral tumor evolution across cancer types. *Nat Genet* **48,**
46   238–244 (2016).
47   30.    Tirosh, I. *et al.* Single-cell RNA-seq supports a developmental hierarchy in
48   human oligodendroglioma. *Nature* (2016). doi:10.1038/nature20123
49   31.    Harvey, J. M., Clark, G. M., Osborne, C. K. & Allred, D. C. Estrogen receptor
50   status by immunohistochemistry is superior to the ligand-binding assay for

predicting response to adjuvant endocrine therapy in breast cancer. *Journal of Clinical Oncology* **17,** 1474–1481 (1999).

32. Buenrostro, J. D. *et al.* Single-cell chromatin accessibility reveals principles of regulatory variation. *Nature* **523,** 486–490 (2015).

33. Cowper-Sal lari, R. *et al.* Breast cancer risk-associated SNPs modulate the affinity of chromatin for FOXA1 and alter gene expression. *Nature Genetics* **44,** 1191–1198 (2012).

34. Cohen, A. J. *et al.* Hotspots of aberrant enhancer activity punctuate the colorectal cancer epigenome. *Nature Communications* **8,** 14400 (2017).

35. Levsky, J. M. & Singer, R. H. Gene expression and the myth of the average cell. *Trends in Cell Biology* **13,** 4–6 (2003).

36. Larson, D. R., Zenklusen, D., Wu, B., Chao, J. A. & Singer, R. H. Real-time observation of transcription initiation and elongation on an endogenous yeast gene. *Science* **332,** 475–478 (2011).

37. Wang, S. *et al.* Target analysis by integration of transcriptome and ChIP-seq data with BETA. *Nature Protocols* **8,** 2502–2515 (2013).

38. Cancer Genome Atlas Network. Comprehensive molecular portraits of human breast tumours. *Nature* **490,** 61–70 (2012).

39. Gyorffy, B. *et al.* An online survival analysis tool to rapidly assess the effect of 22,277 genes on breast cancer prognosis using microarray data of 1,809 patients. *Breast Cancer Res Treat* **123,** 725–731 (2009).

40. Neph, S. *et al.* An expansive human regulatory lexicon encoded in transcription factor footprints. *Nature* **489,** 83–90 (2012).

41. Thurman, R. E. *et al.* The accessible chromatin landscape of the human genome. *Nature* **489,** 75–82 (2012).

42. Paakinaho, V. *et al.* Single-molecule analysis of steroid receptor and cofactor action in living cells. *Nature Communications* **8,** 15896 (2017).

43. Ross-Innes, C. S. *et al.* Differential oestrogen receptor binding is associated with clinical outcome in breast cancer. *Nature* **481,** 389–393 (2012).

44. Kittler, R. *et al.* A comprehensive nuclear receptor network for breast cancer cells. *CellReports* **3,** 538–551 (2013).

45. Carroll, J. S. *et al.* Chromosome-wide mapping of estrogen receptor binding reveals long-range regulation requiring the forkhead protein FoxA1. *Cell* **122,** 33–43 (2005).

46. Beagan, J. A. *et al.* YY1 and CTCF orchestrate a 3D chromatin looping switch during early neural lineage commitment. *Genome Research* **27,** 1139–1152 (2017).

47. Vella, P., Barozzi, I., Cuomo, A., Bonaldi, T. & Pasini, D. Yin Yang 1 extends the Myc-related transcription factors network in embryonic stem cells. *Nucleic Acids Res.* **40,** 3403–3418 (2012).

48. Jeon, Y. & Lee, J. T. YY1 tethers Xist RNA to the inactive X nucleation center. *Cell* **146,** 119–133 (2011).

49. Sigova, A. A. *et al.* Transcription factor trapping by RNA in gene regulatory elements. *Science* **350,** 978–981 (2015).

50. Powe, D. G. *et al.* Investigating AP-2 and YY1 protein expression as a cause of high HER2 gene transcription in breast cancers with discordant HER2 gene amplification. *Breast Cancer Res.* **11,** R90 (2009).

51. Lee, M. H. *et al.* Yin Yang 1 positively regulates BRCA1 and inhibits mammary cancer formation. *Oncogene* **31,** 116–127 (2011).

52. Klymenko, T. *et al.* A Polycomb group protein complex with sequence-specific

DNA-binding and selective methyl-lysine-binding activities. *Genes & Development* **20,** 1110–1122 (2006).

53. Schwalie, P. C. *et al.* Co-binding by YY1 identifies the transcriptionally active, highly conserved set of CTCF-bound regions in primate genomes. *Genome Biology* **14,** R148 (2013).

54. Tang, Z. *et al.* CTCF-Mediated Human 3D Genome Architecture Reveals Chromatin Topology for Transcription. *Cell* **163,** 1611–1627 (2015).

55. Pereira, B. *et al.* The somatic mutation profiles of 2,433 breast cancers refines their genomic and transcriptomic landscapes. *Nature Communications* **7,** 11479 (2016).

56. McLean, C. Y. *et al.* GREAT improves functional interpretation of cis-regulatory regions. *Nature Biotechnology* **28,** 495–501 (2010).

57. Cai, Y. *et al.* YY1 functions with INO80 to activate transcription. *Nature Structural & Molecular Biology* **14,** 872–874 (2007).

58. Hurtado, A., Holmes, K. A., Ross-Innes, C. S., Schmidt, D. & Carroll, J. S. FOXA1 is a key determinant of estrogen receptor function and endocrine response. *Nat Genet* **43,** 27–33 (2010).

59. Vaquero, J., Nguyen Ho-Bouldoires, T. H., Clapéron, A. & Fouassier, L. Role of the PDZ-scaffold protein NHERF1/EBP50 in cancer biology: from signaling regulation to clinical relevance. *Oncogene* (2017). doi:10.1038/onc.2016.462

60. Murtazina, R. *et al.* Tissue-specific regulation of sodium/proton exchanger isoform 3 activity in Na(+)/H(+) exchanger regulatory factor 1 (NHERF1) null mice. cAMP inhibition is differentially dependent on NHERF1 and exchange protein directly activated by cAMP in ileum versus proximal tubule. *Journal of Biological Chemistry* **282,** 25141–25151 (2007).

61. Miller, W. R. *et al.* Gene Expression Profiles Differentiating Between Breast Cancers Clinically Responsive or Resistant to Letrozole. *Journal of Clinical Oncology* **27,** 1382–1387 (2009).

62. Gerlinger, M. *et al.* Genomic architecture and evolution of clear cell renal cell carcinomas defined by multiregion sequencing. *Nature Genetics* **46,** 225–233 (2014).

63. Juric, D. *et al.* Convergent loss of PTEN leads to clinical resistance to a PI(3)Ka inhibitor. *Nature* **518,** 240–244 (2015).

64. Creighton, C. J. *et al.* Development of Resistance to Targeted Therapies Transforms the Clinically Associated Molecular Profile Subtype of Breast Tumor Xenografts. *Cancer Res.* **68,** 7493–7501 (2008).

65. Massarweh, S. *et al.* Tamoxifen Resistance in Breast Tumors Is Driven by Growth Factor Receptor Signaling with Repression of Classic Estrogen Receptor Genomic Function. *Cancer Res.* **68,** 826–833 (2008).

66. Shah, S. P. *et al.* Mutational evolution in a lobular breast tumour profiled at single nucleotide resolution. *Nature* **461,** 809–813 (2009).

67. Bozic, I., Gerold, J. M. & Nowak, M. A. Quantifying Clonal and Subclonal Passenger Mutations in Cancer Evolution. *PLoS Comput Biol* **12,** e1004731 (2016).

68. Mohammed, H. *et al.* Progesterone receptor modulates ERα action in breast cancer. *Nature* **523,** 313–317 (2015).

69. Arduino, S. *et al.* Reduced IL-2 level concentration in patients with breast cancer as a possible risk factor for relapse. *Eur. J. Gynaecol. Oncol.* **17,** 535–537 (1996).

70. García-Tuñón, I. *et al.* Interleukin-2 and its receptor complex (α, β and γ

1   chains) in in situ and infiltrative human breast cancer: an immunohistochemical
2   comparative study. *Breast Cancer Res.* **6,** R1 (2003).
3   71.  Onder, T. T. *et al.* Chromatin-modifying enzymes as modulators of
4        reprogramming. *Nature* **483,** 598–602 (2012).
5   72.  Turner, N. C. *et al.* Palbociclib in Hormone-Receptor–Positive Advanced
6        Breast Cancer. *N Engl J Med* **373,** 209–219 (2015).
7   73.  Schmidt, D. *et al.* ChIP-seq: Using high-throughput sequencing to discover
8        protein–DNA interactions. *Methods* **48,** 240–248 (2009).
9
10

1 **Figure Legends**

2 **Figure 1: Assessment of inter- and intra-tumor epigenetic heterogeneity** A)

3 Main hypothesis of the study. Transcriptional data from bulk tissue represent the

4 average over million cells. Each cell contributes a value from a continuous

5 distribution of potential mRNA molecules. For chromatin data, each cell can only

6 contribute a deterministic value to the bulk signal, generally from two alleles.

7 Therefore, the relative strength of ChIP-seq data is dependent on the number of cells

8 carrying epigenetic signal at discrete loci. H, M and L represent strong, medium and

9 weak signal, respectively. Clonal regulatory regions are commonly shared by BC

10 patients while weak enhancers are more patient specific B) EGR3 mRNA is

11 expressed in MCF7 but not MCF7-F cells. eRNA and Pol-II ChIA-PET show

12 enhancer activity in MC7 but not MCF7-F[16]. CTCF insulated perimeter is shown in

13 yellow. Predicted looping from ChIA-PTE is shown in red. The observed ChIP-qPCR

14 signal for H3K27ac at EGR3 enhancers decrease with increasing number of MCF7-F

15 cells mixed in the sample C) Ranking strategy: H3K27ac signal is normalized at each

16 locus and assigned a ranking index based on relative strength within each single

17 ChIP-seq experiment (1=strongest, 100=weakest, binning on RPKM signal). Binning

18 is repeated for each patient. D) Linear regression shows that clonal enhancers are

19 commonly shared between breast cancer patients. Y axis=Ranking Index, X

20 axis=Sharing index. Sharing index indicate the number of patients sharing the

21 regulatory region. Each dot represents the median RI (all patients) for a single

22 regulatory region. The interpolating lines represent the median RI value and

23 interquartile ranges for regulatory regions with the same SI E) Overlap between BC

24 risk variants and annotated DNA elements F) Variant Set Enrichment analysis

25 indicates that BC-specific but not CRC-specific GWAS risk variants occur more

26 frequently than expected within the enhancers elements identified in our study.

27

28 **Figure 2: Clonal and sub-clonal regulatory regions contain distinct regulatory**

29 **information.** A) Bioinformatic framework of the analyses. H3K27ac calls were split

30 to identify approximate nucleosome-level enrichment (sub-peaks). Sub-peaks data

31 were integrated with ENCODE-derived DHS-seq calls to identify potential sites of TF

32 binding. Individual imputed DHS regions were assigned SI values based on the

1  number of patient sharing the region B) Transcription factor motif analysis of
2  individual bins (SI) followed by unsupervised clustering. RD and RN regions cluster
3  separately in two distinct clades. ERE and YY1 motif are blown up at the bottom C)
4  Clonal enhancers in MCF7 cells (RI<20) are characterized by a higher number of TF
5  footprints, while sub-clonal enhancers (RI>70) have less footprint than expected
6  (O/E=1). Asterisks represent a pValue of <0.001 in a Wilcoxon Signed Rank Test D)
7  Overlap of imputed DHS regions with *in vivo* derived ER$\alpha$ binding sites. The left Y
8  axis indicates cumulative DHS regions. The right Y axes indicate the percentage of
9  overlap based on total DHS in each SI bin E) Distribution plot of *in vivo* derived ER$\alpha$
10  binding sites versus the number of patients in which they were observed[43].

11

12  **Figure 3: YY1 identify a dominant phenotypic clone in ER$\alpha$ BC** A) RIs for the
13  YY1 enhancer within all the individual patients included in the current study. YY1
14  enhancer location with its 3D interactions are shown in the top right inset B) YY1
15  enhancer ranking analysis of available Epigenome Roadmap H3K27ac datasets.
16  Tissues are displayed from the strongest to the weakest YY1 enhancer activity
17  (based on RI). Representative IHC analysis of normal tissues stained with a YY1
18  antibody are shown C) Top left: YY1 expression in ER$\alpha$-positive breast cancer
19  compared to normal breast tissue. Bottom left: Kaplan-Meier analysis of patient
20  outcome using YY1 expression to stratify patients. Right: Kaplan-Meier analysis of
21  patient outcome using YY1 expression. All BC subtypes were analysed separately
22  D) IHC analysis of normal breast tissues highlights YY1 functional subclones in
23  normal breast E) IHC analysis of ER$\alpha$ positive invasive ductal carcinomas identify
24  YY1 positive clones as the dominant clonal population.

25

26  **Figure 4: YY1 marks critical enhancers in breast cancer cells** A) ChIP-seq data
27  from ER$\alpha$-positive MCF7 for YY1 in quiescent or 17ß -estradiol (E2) stimulated cells
28  B) Heatmaps showing global enrichment profiles of several chromatin markers
29  associated with active regulatory regions in MCF7 cells C) Overlap between ER$\alpha$,
30  YY1 and FOXA1 in MCF7 cells. The right panel shows the potential overlap with *in*
31  *vivo-* derived core ER$\alpha$ binding sites D) ER$\alpha$ core binding sites are strongly enriched

1   for YY1 binding in MCF7 cells while patient-specific ERα bindings are generally YY1-

2   free. E) Genes used to classify luminal breast cancer patients are strongly enriched

3   for ERα-YY1 binding sites. Asterisks represent $p<10^{-5}$ in a Fisher's Exact test vs.

4   private ERα F) YY1 depletion leads to transcriptional shut-down of an ERE-driven

5   luciferase reporter. Bars and error bars represent the average of 5 independent

6   experiments with SE. Asterisks represent significance at P<0.001 after ANOVA with

7   Dunnet's correction. G) Silencing YY1 blocks estrogen-induced growth in MCF7 cells

8   H) YY1 depletion leads to growth arrest in AI resistant LTED cells. Proliferation

9   assays were conducted in biological triplicate. Error bars indicate 95% confidence

10  intervals. Asterisks represent significance at P<0.05, 0.01, 0.001 and 0.0001 after 2-

11  way ANOVA with Tukey's post-test I) Overlap of YY1 and ERα binding sites in LTED

12  cell lines J) ERα-YY1 bound enhancers in LTED cells underlie the transcription of

13  genes associated with luminal breast cancer and acquired endocrine therapy

14  resistance K) core ERα-YY1 bound enhancers are strongly enriched near estrogen

15  responsive genes that are not suppressed by Tamoxifen co-treatment.

16

17  **Figure 5: Epigenomic mapping predicts the size of phenotypic clones in**

18  **patients** A) Global Kaplan-Meier analysis summarize univariate analysis for each

19  gene included in the Affymetrix microarray platform. Hazard Ratios are plotted in the

20  X axis B) SLC9A3R1 RNA levels pre- and post- short-term aromatase inhibitor

21  treatment in responder and non-responder patients[61]. Oestrogen-dependent

22  expression of progesterone receptor mRNA is shown as comparison C) Silencing

23  SLC9A3R1 leads to proliferation arrest in response to estrogen stimulation in MCF7

24  and estrogen independent growth in LTED cells. Proliferation assays were

25  conducted in biological triplicate. Error bars indicate 95% confidence intervals.

26  Asterisks represent significance at P<0.05, 0.01, 0.001 and 0.0001 after 2-way

27  ANOVA with Tukey's post-test D) RIs for the SLC9A3R1 enhancer within all the

28  individual patients included in the current study. SLC9A3R1 enhancer location and

29  its 3D interactions are shown in the top right inset E) SLC9A3R1 enhancer ranking

30  analysis of available Epigenome Roadmap H3K27ac datasets. Tissues are displayed

31  from the strongest to the weakest SLC9A3R1 enhancer activity (based on RI).

32  Representative IHC analysis of normal tissues stained with a SLC9A3R1 antibody

are shown. F-G) YY1 and SLC9A3R1 IHC analysis of BC patients profiled using H3K27ac ChIP-seq. Predicted activity (RI) of YY and SLC9A3R1 enhancers is shown on the X axis. The number of cells positively stained for YY1 and SLC9A3R1 protein is indicated on the Y axis. Linear regression R square, confidence intervals and representative staining are also shown.

**Figure 6: Endocrine treatment shapes phenotypic evolution.** A) Theoretical framework of the analysis. The relative size of phenotypic clones can be tracked using enhancer activity (RIs). Phenotypic clones can be positively or negatively selected during BC progression in response to endocrine therapies. B) Expanding or contracting phenotypic clones were defined based on the RI-ratio in primary and metastatic samples ($RI_P/RI_M$). Distribution of RI-ratio identified top candidate enhancers YY1 RI does not change significantly during progression, while SLC9A3R1 RI ranks among the enhancers with stronger increase in activity during progression. Vertical bars represent σ (Standard Deviation) increments from the population median C) Scatterplot of YY1 and SLC9A3R1 enhancer ranking according to patient stage. Bars indicate mean and 95% confidence intervals. Asterisks represent significance at P<0.05 after students two-tail T-Test D) IHC staining for YY1 and SLC9A3R1 positive cells in an independent matched longitudinal cohort of ERα breast cancer patients. All normal and primaries are treatment naïve. All metastatic have received endocrine therapies (Tamoxifen or Aromatase inhibitors). Statistical significance was calculated using a pair-wise, two-tail T-test. Representative images are also shown E) Enhancer and promoter stratification based on frequency of usage in primary and metastatic patients. Percentages were calculated for each regulatory region for each stage (primary and metastatic) and differential was then derived and plotted on the X-axis F) RI indexes for all PE and ME are plotted. As a control, RI for common enhancer (CE) are also plotted. Permutation was used to assess changes in RI in 50 randomly selected sets of CE G) Kaplan-Meier analysis using averaged RNA expression of genes associated with PE or ME regulatory regions. Genes were assigned considering CTCF insulated perimeters E) Pathway analysis for genes associated with PE or ME regulatory regions. Pathways were identified using GREAT and are listed in order of significance (qValue).

39

1

40

**Supplementary Figures**

Supplementary Figure 1. Hematoxylin-Eosin staining to evaluate tumor cellularity was carried out for each sample profiled using ChIP-seq. Only tumors with cellularity above 70% were analyzed.

Supplementary Figure 2. Summary statistics for ChIP-seq analyses. A) The number of individual peaks called using MACS 2.0 are shown for each patient. A q value of 0.01 was used in the peak calling analysis B) Saturation plots. Patients were permutated and the total number of region called after permutation is shown on the Y axis. 80% of total promoters were covered by permutating 4 patients, while similar saturation for enhancers was reached after permutating all 47 samples C) Distribution plots show the frequencies in function of Sharing Index for each regulatory region. Inset show median SI for promoters and enhancers.

Supplementary Figure 3. A) Spiking experiments show that relative enrichment for ERα binding correlates with the number of cells carrying the binding event. MCF7 stimulated or not with estradiol to induce ERα at the specified enhancers were mixed in different proportion with ERα-negative cells with similar genetic background before ERα binding was measured using ChIP-qPCR. Arrows indicate the binding site quantified using ChIP-qPCR B) Linear regression using patient-derived RNA levels and patient-derived SI. The analysis was repeated in three independent cohorts. Normal samples were analysed when available (small insets). The number in each box summarize relative slope for RN and RD elements, the colour of the box indicates the correlation coefficient $R^2$.

Supplementary Figure 4. Transcription factor motif analyses of the entire promoter and enhancer imputed DHS landscapes. Motifs are ranked based on the ratio of observed/expected. Motif were filtered for a q Value of $10^{-4}$

Supplementary Figure 5. YY1 RNA levels from three independent RNA-seq datasets

of normal tissues. Images were obtained using Protein Atlas Tools B) YY1 RNA levels from RNA-seq analysis of cancer cell lines. Images were obtained using Protein Atlas Tools.

Supplementary Figure 6. A) Chromatin landscape at the YY1 enhancer locus in breast cancer cells. The loops were obtained from Pol II ChIA-PET data (high score). The CTCF insulation perimeter was established from CTCF ChIA-PET data (ENCODE). Enhancers SI are shown at the bottom B) Meta-analysis of IHC data from Protein Atlas stained with YY1 antibody. RI for the individual tissues in indicated. Percentage of YY1 positive cells is also listed at the bottom of each image.

Supplementary Figure 7. A) YY1 expression comparing normal tissues and cancer tissues shows that YY1 median expression is significantly stronger in TCGA cancer tissues. Data were generated using TIMER (https://cistrome.shinyapps.io/timer/) (B) YY1 median expression is significantly higher in several breast cancers sub-classes compared to normal tissue. Data were obtained from Oncomine C) IGV snapshot of estradiol induced YY1 ChIP-seq and H3K27ac ChIP-seq from MCF7 cells near transcriptionally inactive genes.

Supplementary Figure 8. A-B) Silencing YY1 is sufficient to abrogate the growth of two independent cell line models of ER$\alpha$ breast cancer. Proliferation assays were conducted in biological triplicate. Error bars indicate 95% confidence intervals. Asterisks represent significance at P<0.05, 0.01, 0.001 and 0.0001 after 2-way ANOVA with Tukey's post-test C) YY1 depletion leads to reduced transcription of common ER$\alpha$ target genes. Each experiment was performed in biological triplicates. Column and bars represent the average and SEM of all experiment. Asterisks represent significance at P<0.05, 0.01, 0.001 and 0.0001 after 1-way ANOVA with Dunnet's post-test.

Supplementary Figure 9. A) SLC9A3R1 RNA expression in BC cell lines sensitive or resistant to endocrine therapies. RNA-seq data were obtained in house[16] B) Meta-analysis of SLC9A3R1 expression in response to ER$\alpha$ depletion C) SLC9A3R1 RNA in additional oestrogen independent BC cell line models. Fold changes are

1   calculated as ratio compared to parental endocrine sensitive BC cells. Oestrogen-

2   dependent expression of progesterone receptor mRNA is shown as comparison.

3   Original microarray codes are shown. Microarray were downloaded from GEO and

4   re-analysed D) SLC9A3R1 transcriptional response to several stimuli is shown. Data

5   were analysed using NURSA Transcriptomime tool

6   (https://www.nursa.org/nursa/transcriptomine/index.jsf;jsessionid=J3hRuy3XjX6HeNr

7   2aekh-rrTsXS-uVUXErsip0wY.nursa3) E) Kaplan Meier survival plots were

8   calculated using three independent large datasets using SLC9A3R1 expression in

9   the primary cancer as a classifier.

10

11  Supplementary Figure10. A) Cell lines from the Protein Atlas initiative were ranked

12  based on SLC9A3R1 expression as profiled by RNA-seq B) SLC9A3R1 expression

13  comparing normal tissues and cancer tissues shows that SLC9A3R1 median

14  expression is significantly stronger in ER$\alpha$-positive Breast cancer patients from

15  TCGA. Data were generated using TIMER (https://cistrome.shinyapps.io/timer/) C)

16  Chromatin landscape at the SLC9A3R1 enhancer locus in breast cancer cells.

17  Looping analysis was conducted using Pol II ChIA-PET data. CTCF insulation

18  perimeter was established from CTCF ChIA-PET data (ENCODE). H3K27ac, YY1,

19  ER$\alpha$, FOXA1 and DHS-seq data were developed in house and deposited online.

20  ER$\alpha$ binding identified as Core are highlighted in the red box D) Expression analysis

21  comparing median RNA expression values for several genes localized near the

22  active SLC9A3R1 putative enhancer. SLC9A3R1 expression is significantly

23  increased in BC tissues compared to normal while the expression of other genes is

24  not affected

25

26  Supplementary Figure11 A) YY1 silencing in MCF7 is sufficient to decrease

27  SLC9A3R1 expression. Each experiment was performed in biological triplicates.

28  Column and bars represent the average and SEM of all experiment. Asterisks

29  represent significance at P<0.05, 0.01, 0.001 and 0.0001 after 1-way ANOVA with

30  Dunnet's post-test B) Silencing SLC9A3R1 is sufficient to abrogate the growth of an

31  endocrine therapy resistant cell line but not the parental, treatment naïve breast

32  cancer cell line. Proliferation assays were conducted in biological triplicate. Error

33  bars indicate 95% confidence intervals. Asterisks represent significance at P<0.05,

1  0.01, 0.001 and 0.0001 after 2-way ANOVA with Tukey's post-test C-D) YY1 and

2  SLC9A3R1 enhancer RIs are shown for all available ENCODE cell lines. MCF7

3  SLC9A3R1 RI is significantly different from the median value calculated on the entire

4  population (without MCF7, green circle). RI and relative RNA levels (RPKM) are

5  shown when available for the same cell type. Asterisks represent significance at

6  $P<0.05$, 0.01, 0.001 and 0.0001 after one sample T-Test.

7

8  Supplementary Figure 12. A) SLC9A3R1 RNA levels from three independent RNA-

9  seq datasets of normal tissues. Images were obtained using Protein Atlas Tools B)

10 SLC9A3R1 enhancer activity was classified based on the relative RI index in each

11 tissue. IHC meta-analysis from the Protein Atlas initiative supports the predicted

12 heterogeneity based on enhancer activity.

13

14 Supplementary Figure 13. A) Protein Atlas sections stained with the indicated

15 antibody were scored and ICH+ positive cells were plotted on the radial axis. Data

16 for each patient are plotted in each corner. Two sections were examined for most

17 patients. Green lines indicate clonal staining.

18

**A** phenotypical heterogeneity and epigenetics

Transcriptome — Genes

Epigenome — Enhancers / Promoters — K27ac

Individual cells

1 2 3 Enhancers

Number of cells in assay

ChIP-seq signal
- H — high % cells — clonal
- M — medium % cells
- L — low % cells
- L — low % cells — sub-clonal

Intra-Tumor

Inter-Patient — Frequency Low → High

**B**

RNA-seq — CTCF insulation

MCF7-F 50rpkm / MCF7 50rpkm

K27ac ChIP-qPCR

PolII ChIA-PET (score 1000)

KIAA1967 — BIN3 — EGR3 — EGR3 eRNA

H3K27ac ChIP-qPCR EGR3 enh

MCF7 / Mix / MCF7-F

Tot input 10⁷ cells — Relative Enrichment (vs. no K27ac site)

100/0  75/25  50/50  25/75  0/100

**C** Bins (by FPKM)

Mets / Primary

Patients — Number of H3K27ac positive regions

highest FPKM — Rank 1 (H) — clonal
2, 3
Rank 100 (L) — lowest FPKM — sub-clonal

Rank 1 (H) : Rank 100 (L) — Ranking Index (RI)

**D**

Rank Index (Patients H3K27ac)
clonal 1 / 50 / sub-clonal 100

Promoters — R²=0.98 SLOPE=1.57

Enhancers — R²=0.97 SLOPE=1.22

Sharing Index (SI)
1 2 3 — Regulatory Noise — 10 — 20 — 30 — 40 — 45 Regulatory Drivers
Private → Some patients → Most patients

**E** BC Risk SNPs Overlap

UTR
TSS
Introns
Promoters
Exons
SI=>21
SI=1-21
SI=1

**F**

BC risk SPNs / CRC risk SPNs

UTR
TSS
Introns
Promo.
Exons
SI=>21
SI=1-21
SI=1

Enrichment Score (VSE)

**A**

H3K27ac (patients) | DHS-seq (129 cell lines) | Patient specific call | SI annotation

peak splits

potential TF B.S.

merge

#1
#2
#..
#n

E  P

#1-6, 18, 24
#3, 19, 40

Motif analysis

**B**

Imputed TF analysis (SI specific)
promoters

Observed/Expected qVal<0.05
0  1  2  >3

YY1
JUND
PBX1
PKNOX1
GATA3
ZNF143
ELF5
SPDEF
GATA3
HOXB13
FOXM1
FOXA1

MEIS1

Sharing Index
1  10  20  30  40  >40

RN          RD

Imputed TF analysis (SI specific)
enhancers

Observed/Expected qVal<0.05
0  1  2  >3

GATA3
ERE

YY1

JUN/AP
FRA1
IRF2
IRF1
SPDEF
FOXA1      CTCF
FOXM1
IRF1

RD          RN

YY1

ERE

**C**

DHS Footprints

Footprint Enrichment
0      1      4

MCF7 Ranking Index
1
10
20
30
40
50
60
70
80
90
100

***

***

O/E<1      O/E>1

**D**

*in vivo* ERα ChIP-seq (all)

Cumulative n° DHS /K27ac. regions
10^6
10^5
10^4
10^3

Sharing Index
1          45

100%
10%
1%

% overlap with ERα Binding

10^6
10^5
10^4
10^3

Sharing Index
1        45

100%
10%
1%

**E**

Total *in vivo* ERα binding

Most patients

Good Outcome
Poor Outcome

ERα core

ERα single patients

9

ERα Sharing Index x 9

1

Private

100%      80%          20%        0%
Percentage of total ERα bound regions (in vivo)

## A. BCa Patients: YY1 enhancer A
### Ranked H3K27ac ChIP-seq signal

○ Primary
○ Metastatic

YY1 enhancer A (SI=41)
CTCF insulation
enhancer A
DEGS2    YY1    SLC25A29
25kb

H / M / L
R.I.
1, 20, 80, 100

clonal sub-clonal

patient number
1 — 10 — 20 — 30 — 40

## B. Normal Tissues: YY1 enhancer A
### Epigenetic Roadmap Normal Tissues. H3K27ac ChIP-seq

YY1 enhancer A (SI=41)
CTCF insulation
DEGS2    YY1
25kb

H / M / L
R.I.

Ovary, Esophagus, Spleen, Large Intestine, Aorta, Heart Left Ventricle, Mesoderm, Lung, Heart Right Vent., Placenta, Sigmoid Colon, Small Intestine, CD8, T cells, CD4, helper T cells, Thymus, Psoas Muscle, Right Atrium, Stomach, NK cells, Adipose, Myeloid CD34, Urinary Bladder, Adrenal Gland, T-Cells, CD4, memory T cells, CD14-positive monocytes, H1-Stem Cells, Neuronal Stem Cells, B cells, Trophoblast, HUES64, Muscle of the trunk, I9 cells, IMR-90, iPS-18a, iPS-20b, Mesenchymal cells

H — Ovary
H — Lung
M — Stomach
M — Adrenal
L — Skeletal Muscle

## C.
### METABRIC YY1 RNA

Norm. Expression
12 / 7
pV=1.9$^{-24}$

Normal n-=144
ER+ BCa n=1486

### TCGA O.S.

Lum. A — 121 pt HR=3 **
Lum. B — 670 pt HR=1.4
HER2 — 67 pt HR=0.33
TNBC — 209 pt HR=0.55

% Survival

2 4 6 years

### METABRIC ER+ O.S.

% Survival
100 / 0

YY1$^{high}$=15.2 yrs.
YY1$^{low}$=18.7 yrs.

pV=0.0001
HR=1.43
****

10 — 20 years

## D.
Normal Lobule    Normal Duct

YY1+    YY1−
YY1+    YY1−
YY1−    YY1+

## E.
ERα IDC    Cancer    Duct
ERα IDC    Cancer    Duct    Cancer

**A** MCF7 YY1 ChIP-seq
(consensus 2 biological replicates)

GREAT ETOH

RNA binding · 1e-28
mRNA processing · 1e-37
mRNA splicing · 1e-17

18.2% Promoters
81.8% Distal
+ETOH
+E2

164   2059   45970

GREAT E2

ER-nucleus signaling pathway · 1e-63
Histone methyltransferase complex · 1e-24
Genes UP in lum. vs. mesen. BCa cell lines · 1e-260
Genes UP in lum. vs. basal BCa cell lines · 1e-184

Mitochondrial genes · 1e-32
Cell cycle specific genes · 1e-24

YY1 · 1e-1167 · 51% motif
ELK4 · 1e-45 · 32% motif

MOTIF ETOH

YY1 · 1e-751 · 10% motif
CTCF · 1e-752 · 13% motif

MOTIF E2

**B** G2M H3K27ac G1   YY1 ETOH   YY1 E2   DHS   ERαE2 FOXA1

YY1 2223
YY1 E2 45970

-5kb 5kb

ETOH  YY1 E2  H3K27ac
FOXA1

ETOH  YY1 E2  H3K27ac
CCDC117 XBP1

ETOH  YY1 E2  H3K27ac
B   A C   YY1

**C** MCF7 YY1 and ERα

ERa +E2
YY1 +E2
7.4  1  37.8
3  2
9.7
49.5
FOXA1

F
Y
E
FY
FE
YE
FYE

% overlap with core ERα (484)

0   10

**D** In vivo ERα and YY1
P val 1 10⁻⁵⁰

303/484
36.9K/215.9K

% overlap with E2-YY1
core ERα   Private ERα

**E** TCGA luminal signature

6 ****  ****

Enrichment over null signature

private ERα   core ERα   core ERα+YY1

**F**
EtOH   E2
YY1
GAPDH
siRNA   -   +a  +b   -  +a  +b

ERE-Luc Reporter
EtOH   E2
0.6  ***
0.4
0.2

R.U. Luciferase
mock   siYY1ª   siYY1ᵇ

**G** siYY1 proliferation
MCF7
0.4
0.3 ****
0.2 *
0.1

Relative Units
0   3   5   7   days

○ mock
● mock+E2
○ siYY1ª
● siYY1ª+E2
○ siYY1ᵇ
○ siYY1ᵇ+E2

**H** siYY1 proliferation
LTED
0.4
0.3 ****
0.2 ****
0.1

0   3   5   7   days

○ mock
○ siYY1
○ siYY1ᵇ

**I** LTED ChIP-seq

52.1   5   20
YY1   ERα

**J** GREAT pathways
YY1-ERα bound LTED enhancers

200

-Log(Q)

50

UP Luminal vs. Basal BCa
SMAD2 target in HaCat
Genes bound by ER and up by E2
Down in Basal BCa
Acquired ETR in ESR1 and ERBB2 BCa

**K** TAM resistant estrogen dependent gene (84/723)

TSS   exon   3UTR
-20Kb   region examined for each gene   3' end

20

15 ****

10 ****

5
** ** *

Enrichment vs. Null List

IN VIVO private ERα
MCF7 YY1-E2
MCF7 ERα and YY1-E2
CORE ERα
CORE ERα and YY1

CXXC5
SLC9A3R1
FAM102A
FKBP4
MYB
KRT13
CA12
CELSR2
RAPGEFL1
MAG

**A** 724 ER+ BCa treated endocrine therapy

H.R. (relapse)
slower ← → faster
○ all genes
● SLC9A3R1

**B** Neo-Adjuvant AI treatment

SLC9A3R1 mRNA

PGR mRNA

pre- post-
Responder 37 patient

pre- post-
Non-Responder 14 patients

**C** siSLC9A3R1 proliferation

MCF7
○ mock
● mock+E2
○ siSLC1
◐ siSLC1+E
○ siSLC1b
○ siSLC1b +E2
****

LTED
days
○ mock
◐ siSLCa
○ siSLCb
****

days

**D** BCa Patients: SLC9A3R1 enhancer
Ranked H3K27ac ChIP-seq signal

○ Primary
○ Metastatic

SLC9A3R1 enhancer (SI=34)
CTCF insulation
RAB37   SLC9A3R1   NAT9
14kb

H, M, L  R.I.  clonal sub-clonal

patient number

**E** SLC9A3R1 enhancer heterogeneity
Epigenetic Roadmap Normal Tissues. H3K27ac ChIP-seq

SLC9A3R1 enhancer (SI=34)
CTCF insulation
RAB37   SLC9A3R1   NAT9
14kb

H, M, L  R.I.

Esophagus, Trophoblast, Small Intestine, Large Intestine, Neuronal Stem Cell, Mesoderm, Thymus, Placenta, Stomach, Urinary Bladder, Heart Left Ventricle, Right Cardiac Atrium, H1-Stem Cells, Adrenal Gland, CD8, T cells, Sigmoid Colon, T-Cells, HUES64, NK cells, Myeloid progenitor CD34, Spleen, CD4 Helper T cells, Mesenchymal cells, CD4, memory T cells, Muscle of the trunk, Psoas Muscle, Adipose, Aorta, B cells, CD14-positive monocytes, H9 cells, Heart Right Ventricle, IMR-90, PS-18a, Lung, Ovary

H  Duodenum
M  Stomach (upper)
M  Adrenal Gland
M  Spleen
L  Adipose
L  Lung

**F** BCa Patients: YY1 protein

% IHC + cells

Mets
Primary
clonal  sub-clonal

R.I.
20    80

**G** BCa Patients: SLC9A3R1 protein

% IHC + cells

Mets
R²=0.74
HPA9672
Primary
clonal  sub-clonal

R.I.
20    80

HPA9672
P14 RI=20
P11 RI=54

HPA1119  P17 RI=8    P33 RI=10.5    P15 RI=13    M003 RI=13    P29 RI=100 POS NEG    P26 RI=76 NEG POS NEG

**A** Tracking phenotypic clones through BC progression

Surgery  Relapse

Primary

RI changes

Metastasis

Candidate Enhancer 1 R.I.

100% H

M

L

0%

Endocrine Therapy

Candidate Enhancer 1 % cells/sample

100% H

L  L

0%

RI increase
RI decrease

**B**

RI increase  RI decrease

Frequency

Yy1
(83293/88935)
-107σ

SLC9A3R1
(54/88935)
+3.86 σ

Log2 ΔRI (Ri_P/Ri_M)

-2  2

**C**

YY1 enhancer A
(discovery dataset)

clonal sub-clonal

R.I

20

80

primaries  mets

SLC9A3R1 enhancer A
(discovery dataset)

P=0.045

clonal sub-clonal

R.I

20

80

primaries  mets

**D**

YY1 IHC Matched P-M
(validation dataset)

clonal sub-clonal

% IHC + cells

100

patient 1  +  −

patient 2

primary  metastatic

Endocrine Therapy

SLC9A3R1 IHC Matched P-M
(validation dataset)

clonal sub-clonal

% IHC + cells

100

ERαloss

P<0.0001

patient 1  +  −  −  +  +

patient 2

primary  metastatic

Endocrine Therapy

**E**

Primary Specific  differential % distribution  Metastas Specific

Promoters  239  20

PE  CE  ME

Enhancers  324 (1‰)  301 (1‰ and P<0.05)

40%  20%  0%  20%  40%

**F**

PE  Rank Index  1  ***

CE  Rank Index  1

ME  Rank Index  1  ***

100  100  100

CE

2  5  10  15  20

PE
ME

-Log(pV) (50 rand, 320 CE sets)

**G** KM with Enhancers-associated genes
METABRIC ER+ (1427)

100  P^high=17.3 yrs.
P^low=12.5 yrs.

% Survival

Transcripts associated with PE

****

10  20

100  M^high=15.3 yrs.
M^low=17.3 yrs.

% Survival

Transcripts associated with ME

*

10  20  years

**H** Enhancer-associated pathways

-Log(Qval)  10  2

Primary

Up in T24 upon PDT stress
Down in AIL T vs. T cells
Down in multiple epithelial tumors
Up in MCF7 upon NRGα
Bound by FOXP3 in hybridoma
UP in MM1S upon adophostin
Up in HeLa upon EGF
Up in CaLu-6 upon TNF
Bound by FOXP3 in hybridoma upon PMA
Down in SEN D upon si ELK3
Down in OV90 upon invasion inhibition
Up in MCF7 upon EGF
IL-2 mediated signaling
Pathological Neovascularization
Abnormal cell proliferation

-Log(Qval)  15  5

Metastasis

Up in ERBB2 BCa
Genes in amplicon 20q12-q13 in BCa
Up in acquired ET ER no ERBB2
Luminal BCa
Up in Luminal BCa
Up in MCF7/DU145 upon BRCa1
Up in MCF7 xenograft TAM resistant
PPARG/RXRA adipocite
Up Luminal BCa cells
EGF target in 184a1 (mammary)
Down in Basal BCa
Up in acquired ET ER and ERBB2
FOXA1 network