

Research article

Open Access

## Gene expression signature of estrogen receptor $\alpha$ status in breast cancer

Martín C Abba<sup>1</sup>, Yuhui Hu<sup>1</sup>, Hongxia Sun<sup>1</sup>, Jeffrey A Drake<sup>1</sup>, Sally Gaddis<sup>1</sup>, Keith Baggerly<sup>2</sup>, Aysegul Sahin<sup>3</sup> and C Marcelo Aldaz\*<sup>1</sup>

Address: <sup>1</sup>Department of Carcinogenesis, The University of Texas M.D. Anderson Cancer Center, Science Park-Research Division, Smithville, Texas, USA, <sup>2</sup>Department of Biostatistics, The University of Texas M.D. Anderson Cancer Center, Houston, Texas, USA and <sup>3</sup>Department of Pathology, The University of Texas M.D. Anderson Cancer Center, Houston, Texas, USA

Email: Martín C Abba - mabba777@hotmail.com; Yuhui Hu - yhhu@mdanderson.org; Hongxia Sun - hcsun@mdanderson.org; Jeffrey A Drake - jadrake@mdanderson.org; Sally Gaddis - sgaddis@mdanderson.org; Keith Baggerly - kabagg@mdanderson.org; Aysegul Sahin - asahin@mdanderson.org; C Marcelo Aldaz\* - maldaz@odin.mdacc.tmc.edu

\* Corresponding author

Published: 11 March 2005

Received: 03 November 2004

BMC Genomics 2005, 6:37 doi:10.1186/1471-2164-6-37

Accepted: 11 March 2005

This article is available from: <http://www.biomedcentral.com/1471-2164/6/37>

© 2005 Abba et al; licensee BioMed Central Ltd.

This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/2.0>), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

### Abstract

**Background:** Estrogens are known to regulate the proliferation of breast cancer cells and to modify their phenotypic properties. Identification of estrogen-regulated genes in human breast tumors is an essential step toward understanding the molecular mechanisms of estrogen action in cancer. To this end we generated and compared the Serial Analysis of Gene Expression (SAGE) profiles of 26 human breast carcinomas based on their estrogen receptor  $\alpha$  (ER) status. Thus, producing a breast cancer SAGE database of almost 2.5 million tags, representing over 50,000 transcripts.

**Results:** We identified 520 transcripts differentially expressed between ER $\alpha$ -positive (+) and ER $\alpha$ -negative (-) primary breast tumors (Fold change  $\geq 2$ ;  $p < 0.05$ ). Furthermore, we identified 220 high-affinity *Estrogen Responsive Elements* (EREs) distributed on the promoter regions of 163 out of the 473 up-modulated genes in ER $\alpha$  (+) breast tumors. In brief, we observed predominantly up-regulation of *cell growth* related genes, *DNA binding* and *transcription factor activity* related genes based on Gene Ontology (GO) biological functional annotation. GO terms over-representation analysis showed a statistically significant enrichment of various transcript families including: *metal ion binding* related transcripts ( $p = 0.011$ ), *calcium ion binding* related transcripts ( $p = 0.033$ ) and *steroid hormone receptor activity* related transcripts ( $p = 0.031$ ). SAGE data associated with ER $\alpha$  status was compared with reported information from breast cancer DNA microarrays studies. A significant proportion of ER $\alpha$  associated gene expression changes was validated by this cross-platform comparison. However, our SAGE study also identified novel sets of genes as highly expressed in ER $\alpha$  (+) invasive breast tumors not previously reported. These observations were further validated in an independent set of human breast tumors by means of real time RT-PCR.

**Conclusion:** The integration of the breast cancer comparative transcriptome analysis based on ER $\alpha$  status coupled to the genome-wide identification of high-affinity EREs and GO over-representation analysis, provide useful information for validation and discovery of signaling networks related to estrogen response in this malignancy.

## Background

Estrogen plays essential roles in the development, growth control and differentiation of the normal mammary gland. However, it is well documented that endogenous estrogens are powerful mitogens critical for the initiation and progression of human breast and gynecological cancers [1]. This cell proliferation signal is mediated by the estrogen receptors (ER), members of the nuclear receptor family that function both as signal transducers and transcription factors to modulate expression of target genes [2]. There are two main subtypes of estrogen receptors: ER $\alpha$  and ER $\beta$  that generally can form homo- and heterodimers before binding to DNA. Although the DNA binding domains of these receptors are very similar, the overall degree of homology is low [3].

Transcriptional regulation of target genes in response to 17 $\beta$ -estradiol (E<sub>2</sub>) is mediated by two main mechanisms. In one, the E<sub>2</sub>-ER complex binds to a specific DNA sequence called the estrogen response element (ERE), this receptor-ligand DNA bounded complex interacts with co-regulatory proteins, promoting chromatin remodeling and bridging with the general gene transcription machinery thus resulting in transcription initiation [4]. Alternatively, the ligand-ER complex can interact with other DNA-bound transcription factors that in turn bind DNA sequences (e.g. via AP1, SP1 complexes) [5,6]. ER $\alpha$  and ER $\beta$  have different affinities for different response elements and exhibit distinct transcriptional properties. Additionally, E<sub>2</sub> also exerts rapid, non-genomic effects attributed to cell membrane-initiated signaling [7].

Approximately two-thirds of all breast cancers are ER $\alpha$  (+) at the time of diagnosis and expression of this receptor is determinant of a tumor phenotype that is associated with hormone-responsiveness. Patients with tumors that express ER $\alpha$  have a longer disease-free interval and overall survival than patients with tumors that lack ER $\alpha$  expression [8]. However, the association between ER $\alpha$  expression and hormonal responsiveness is not perfect: approximately 30% of ER $\alpha$ -positive tumors are not hormone-responsive while 5–15% of ER $\alpha$ -negative tumors respond to hormonal therapy [9]. The molecular basis for the association between ER $\alpha$  expression, hormonal responsiveness and breast cancer prognosis remains unclear.

Several studies have been carried out using cDNA and oligonucleotide microarrays identifying breast cancer subclasses possessing distinct biological and clinical properties [10-13]. Among the distinctions made to date, the clearest separation was observed between ER $\alpha$  (+) and ER $\alpha$  (-) tumors [10-15]. It has been suggested that there are sets of genes expressed in association with ER $\alpha$  that could play an important role in determining the hor-

mone-responsive breast cancer phenotype [16]. ER $\alpha$  is obviously likely to be important for the E<sub>2</sub> induced proliferative response predominantly via the regulation of estradiol-responsive genes. Nevertheless, the expression of additional subsets of genes not necessarily directly regulated by estrogen may also be fundamental in defining the breast cancer hormone-responsive phenotype.

To further elucidate the molecular basis of estrogen-dependent breast carcinogenesis, we here report a comparative transcriptome profiling of invasive breast tumors based on ER $\alpha$  status obtained by SAGE. The SAGE method provides a statistical description of the mRNA population present in a cell without prior selection of the genes to be studied, and this constitutes a major advantage [17]. The breast cancer SAGE comparative analysis was combined with promoter sequence analysis of genes of interest using high-throughput methods of high-affinity ERE identification. In order to have an even more comprehensive picture we also performed a cross-platform comparison between SAGE and DNA microarray studies.

## Results and discussion

### Biomarkers of ER $\alpha$ status in breast carcinomas

The primary goal of our study was to identify the most commonly deregulated genes in invasive breast carcinomas related to ER $\alpha$  status. To this end SAGE data was obtained from a set of primary breast carcinomas. Thus, a breast cancer SAGE database of almost 2.5 million tags was analyzed, representing over 50,000 tag species. We performed a comprehensive evaluation and comparison of gene expression profiles using a recently developed supervised method [18], to identify the most representative differentially expressed transcripts between tumors groups, *i.e.* ER $\alpha$  (+) vs. ER $\alpha$  (-) breast tumors.

This statistical analysis revealed 520 genes differentially expressed (Fold change  $\geq 2$ ;  $p < 0.05$ ) between ER $\alpha$  (+) and ER $\alpha$  (-) primary breast carcinomas (see additional data file 1). Among the 520 transcripts, 473 were up-modulated and 47 were down-modulated transcripts in ER $\alpha$  (+) tumors.

The most commonly over-expressed transcripts in ER $\alpha$  (+) tumors were: *trefoil factor 1 (TFF1/pS2)*, *synaptotagmin-like 4 (SYTL4)*, *regulating synaptic membrane exocytosis 4 (RIMS4)*, *dual specificity phosphatase 4 (DUSP4)*, *chromosome 1 open reading frame 34 (C1orf34)*, *neudin homolog (NDN)*, *n-acetyltransferase 1 (NAT1)* and *caspase recruitment domain family 10 (CARD10)* (Table 1 and additional data file 1).

To validate novel ER $\alpha$  associated genes detected by SAGE not reported in other studies, we performed Real Time RT-PCR analysis of representative transcripts in an

**Table 1: Most highly up-modulated transcripts in ER $\alpha$  (+) breast carcinomas identified by SAGE.**

Gene name	Tag	Locus Link	Fold change (p value)	Frequency#
<b>Cell proliferation related</b>				
<i>TFF1</i> * (trefoil factor 1)	CTGGCCCTCG	7031	51.4 (0.0016)	15/18 (83%)
<i>DUSP4</i> (dual specificity phosphatase 4)	CGGGCAGAAA	1846	14.7 (0.0016)	14/18 (78%)
<i>NDN</i> * (necdin homolog)	ACCTTGCTGG	4692	13.3 (0.0026)	11/18 (61%)
<i>HDGFRP3</i> (hepatoma-derived growth factor)	TGTAAAGTTT	50810	9.8 (0.0019)	12/18 (67%)
<i>TSPAN1</i> * (tetraspan 1)	GGAAGTGTGA	10103	9.5 (0.0017)	15/18 (83%)
<i>SEP6</i> (septin 6)	TCAATTTTCA	23157	7.6 (0.0044)	12/18 (67%)
<i>DHX34</i> * (DEAH box polypeptide 34)	GTTGCTCACT	9704	7.1 (0.0129)	9/18 (50%)
<b>Apoptosis related</b>				
<i>CARD10</i> * (caspase recruitment domain family)	AGAATGTACG	29775	11.1 (0.0030)	15/18 (83%)
<b>Signal transduction related</b>				
<i>SYTL4</i> * (synaptotagmin-like 4)	TATGTGTGCT	94121	28.0 (0.0003)	15/18 (83%)
<i>ECM1</i> * (extracellular matrix protein 1)	ACTGCCCGCT	1893	10.1 (0.0175)	13/18 (72%)
<i>LEPR</i> * (leptin receptor)	AAAGTTTGAG	3953	9.8 (0.0302)	10/18 (55%)
<i>PTGES</i> (prostaglandin E synthase)	TGAGTCCCTG	9536	8.0 (0.0168)	8/18 (44%)
<i>SCUBE2</i> (signal peptide, CUB domain EGF-like 2)	TCAGCACAGT	57758	7.5 (0.0024)	14/18 (78%)
<i>ADORA2A</i> * (adenosine A2a receptor)	TGCTGAGTAG	135	7.1 (0.0460)	11/18 (61%)
<i>ITGBL1</i> (integrin beta-like 1)	CATATTCACA	9358	7.1 (0.0159)	8/18 (44%)
<b>Regulation of transcription related</b>				
<i>ESR1</i> (estrogen receptor 1)	AGCAGGTGCC	2099	9.8 (0.0000)	18/18 (100%)
<i>TCEAL1</i> (transcriptional elongation factor A)	AAAGATGTAC	9338	9.8 (0.0014)	13/18 (72%)
<i>ZNF14</i> (zinc finger protein 14)	TAAACAGCCC	7561	8.4 (0.0023)	13/18 (72%)
<i>ZNF38</i> * (zinc finger protein 38)	CCAGCATTAC	7589	7.6 (0.0051)	10/18 (55%)
<i>HIF1AN</i> * (hypoxia-inducible factor 1 $\alpha$ subunit inhibitor)	CCTGAGTGCG	55662	7.1 (0.0094)	10/18 (55%)
<i>HOXC13</i> (homeo box C13)	TTTTTAAAAT	3229	7.1 (0.0157)	9/18 (50%)
<b>Cytoskeleton</b>				
<i>MAPT</i> (microtubule-associated protein tau)	GTAGACTCGC	4137	9.8 (0.0085)	9/18 (50%)
<i>MYLIP</i> (myosin regulatory light chain interacting)	TTTTCCACTC	29116	9.3 (0.0036)	11/18 (61%)
<b>Metabolism and Miscellaneous</b>				
<i>RIMS4</i> (regulating synaptic membrane exocytosis)	TTGAAATTAA	140730	24.9 (0.0378)	8/18 (44%)
<i>NATI</i> (N-acetyltransferase 1)	TATCTTCTGT	9	11.7 (0.0385)	15/18 (83%)
<i>ATP6V1B1</i> * (ATPase, H <sup>+</sup> transporting)	CCTCCCCCTC	525	10.7 (0.0111)	10/18 (55%)
<i>JDPI</i> (J domain containing protein 1)	TCTGTGAATT	56521	10.0 (0.0035)	12/18 (67%)
<i>CHST11</i> (carbohydrate sulfotransferase 11)	AACCTTCCTC	50515	9.8 (0.0009)	13/18 (72%)
<i>CILP</i> (nucleotide pyrophosphohydrolase)	GTTTTGCCCA	8483	9.3 (0.0054)	14/18 (78%)
<i>ABCA3</i> (ATP-binding cassette sub-family A)	GTAGTCACCG	21	8.9 (0.0149)	10/18 (55%)
<i>SEC14L2</i>	GGAAGGCGGC	23541	8.7 (0.0487)	9/18 (50%)
<i>ANXA9</i> * (annexin A9)	ACATCCGAGG	8416	8.4 (0.0145)	10/18 (55%)
<i>KCTD3</i> (K channel tetramerisation domain 3)	ATAATTAAT	51133	8.4 (0.0001)	17/18 (94%)
<i>SFRS7</i> (splicing factor)	TAGCTAATAT	6432	8.0 (0.0031)	12/18 (67%)
<i>SNRPA</i> * (small nuclear ribonucleoprot. polypep. A)	AAGATCTCCT	6626	7.6 (0.0009)	15/18 (83%)
<i>NNMT</i> (nicotinamide N-methyltransferase)	CCTGCAATTC	4837	7.6 (0.0120)	10/18 (55%)
<i>SLC1A4</i> (solute carrier family 1 member 4)	GACTCACAGG	6509	7.6 (0.0254)	9/18 (50%)
<i>TIPARP</i> (TCDD-inducible polymerase)	AAATGGCCAA	25976	7.6 (0.0051)	10/18 (55%)
<i>SLC7A2</i> (solute carrier family 7 member 2)	CACTGACAGC	6542	7.3 (0.0190)	11/18 (61%)
<i>GA</i> * (liver mitochondrial glutaminase)	CTGCTGCTAC	27165	7.1 (0.0126)	9/18 (50%)
<b>Function unknown</b>				
<i>C1orf34</i>	AGGATGTACA	22996	13.3 (0.0025)	14/18 (78%)
<i>SMILE</i> (hypothetical protein FLJ90492)	TAGAGAGTTT	160418	11.1 (0.0004)	15/18 (83%)

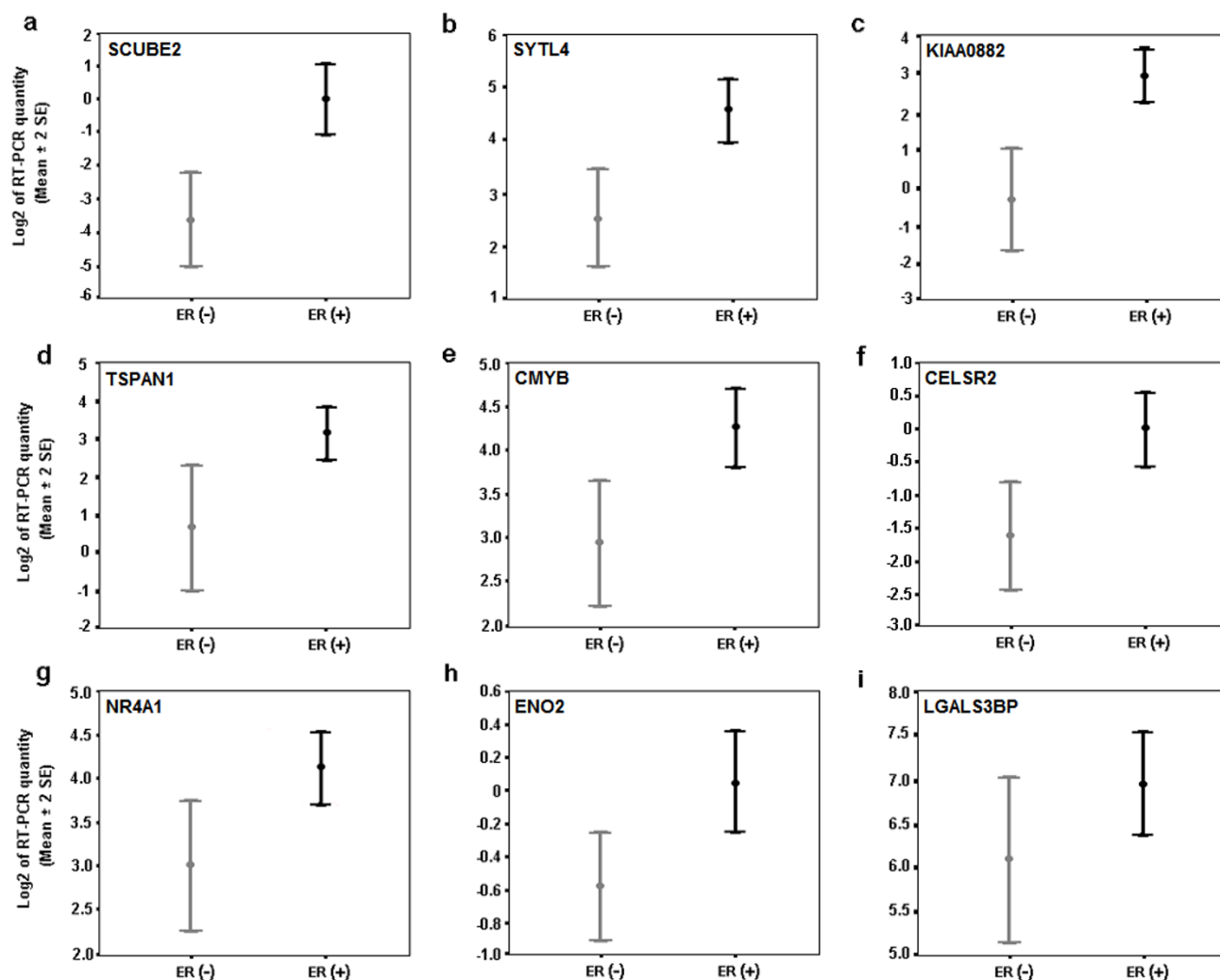
**Table 1: Most highly up-modulated transcripts in ERα (+) breast carcinomas identified by SAGE. (Continued)**

<i>RHBDL4</i> (rhomboïd, veinlet-like 4)	TTGTTTCTAA	162494	10.7 (0.0099)	9/18 (50%)
<i>KIAA0882</i>	GTCTCATTTTC	23158	10.1 (0.0007)	18/18 (100%)
<i>C20orf103*</i>	TTTAGTGATT	24141	9.3 (0.0277)	10/18 (55%)
<i>FLJ33387</i>	GCAGGGAGAG	161145	9.3 (0.0118)	10/18 (55%)
<i>TRALPUSH</i>	GTTTCCAGAG	116931	8.9 (0.0458)	9/18 (50%)
<i>KIAA0980*</i>	TGGTGCTTCC	22981	7.6 (0.0096)	11/18 (61%)
<i>C10orf32</i>	AGTCTGTTGT	119032	7.3 (0.0002)	15/18 (83%)
<i>FLJ13611</i>	TAATCACACT	80006	7.1 (0.0069)	10/18 (55%)

\* Genes with known or putative high-affinity EREs mapping in the vicinity of the TSS.

# Transcripts tags changing > 2-fold when compared with the average expression of ER (-) tumors in at least 8 of 18 (44%) ERα (+) invasive carcinomas SAGE libraries.

For the whole list of ERα associated transcripts see additional data file 1.



**Figure 1**

**Real time RT-PCR validation of nine over-expressed genes in 36 invasive breast carcinomas. a) *SCUBE2* (p = 0.0001); b) *SYTL4* (p = 0.0005); c) *KIAA0882* (p = 0.0005); d) *TSPAN1* (p = 0.001); e) *CMYB* (p = 0.002); f) *CELSR2* (p = 0.011); g) *NR4A1* (p = 0.029); h) *ENO2* (p = 0.033); i) *LGALS3BP* (p = 0.079).** Mean ± 2 Standard Error based on Log2 transformation of real time RT-PCR values of the assayed gene relative to 18S rRNA used as normalizing control.

independent set of 36 invasive ductal breast carcinomas. In agreement with our SAGE analysis, we detected statistical differences in the over-expression of 8 out of 9 evaluated transcripts in ER $\alpha$  (+) breast tumors including: *signal peptide CUB domain EGF-like 2 (SCUBE2)* ( $p = 0.0001$ ), *SYTL4* ( $p = 0.0005$ ), *KIAA0882 protein* ( $p = 0.0005$ ), *tetraspan 1 (TSPAN1)* ( $p = 0.001$ ), *myeloblastosis viral oncogene homolog (C-MYB)* ( $p = 0.002$ ), *epidermal growth factor-like 2 (CELSR2)* ( $p = 0.011$ ), *nuclear receptor subfamily 4 (NR4A1)* ( $p = 0.029$ ), and *enolase 2 (ENO2)* ( $p = 0.033$ ) (Figure 1). A trend of borderline significance was detected for the *lectin galactoside-binding protein (LGALS3BP)* ( $p = 0.079$ ) transcript (Figure 1).

*SCUBE2* (also known as *EGF-like 2* or *CEGP1*) encodes a secreted and cell-surface protein containing EGF and CUB domains that defines a novel gene family [19]. The epidermal growth factor (EGF) motif is found in many extracellular proteins that play an important role during development, functioning as secreted growth factors, transmembrane receptors, signaling molecules, and important components of the extracellular matrix. The CUB domain is found in several proteins implicated in the regulation of extracellular process such as cell-cell communication and adhesion [20]. Expression of *SCUBE2* has been detected in vascular endothelium and may play important roles in development, inflammation and perhaps carcinogenesis [19].

The *CELSR2* gene (also known *EGFL2*) encodes a protein member of the nonclassic-type cadherins (flamingo subfamily). These 7-pass transmembrane proteins have nine cadherin domains, seven-epidermal growth factor-like repeats and two laminin A G-type repeats [21]. It is postulated that these proteins are receptors involved in cell adhesion and receptor-ligand interactions [21] playing a role in developmental processes and cell growth/ maintenance in epithelial and neuronal cells [22,23].

*SYTL4* (also known as *granuphilin-a* or *SLP4*) contains an N-terminal Slp homology domain (SHD) than can specifically and directly bind the GTP-bound form of Rab27A, a small GTP-binding protein involved in granule exocytosis in cytotoxic T lymphocytes [24]. We determined that over-expression of *SYTL4* is associated with ER $\alpha$  (+) tumors (Figure 1b). However, the potential role of this gene in breast carcinogenesis remains unknown.

*ENO2* (also known as *NSE/neuron-specific gamma enolase*) encodes one of three enolase isoenzymes found in mammals. This isoenzyme was described to be expressed in cells of neuronal origin. Interestingly, in a recent report Hao *et al.* (2004) showed high expression of *ENO2* transcripts in breast cancer lymph node metastases when compared with primary breast tumors [25].

The *TSPAN1* gene (also known as *tetraspanin* or *NET1*) encodes a cell-surface protein member of the transmembrane 4 superfamily (*TM4SF*), involved in the regulation of cell development, activation, growth and motility. A number of tetraspanins were described as tumor-specific antigens, and it was suggested that the function of some *TM4SF* proteins may be particularly relevant to tumor cell metastasis [26]. Sugiura and Berditchevski (1999) observed that *TSPAN1* protein complexes may control the invasive migration of tumor cells and contribute to ECM-induced production of MMP2 in breast cancer cell line [27].

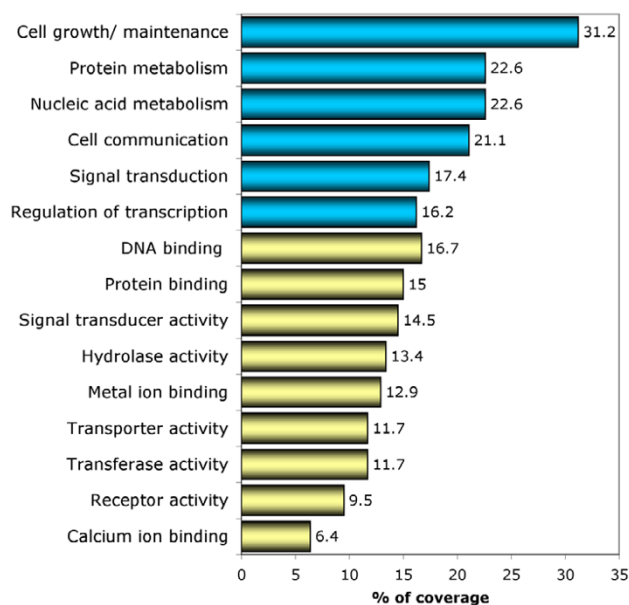
*NR4A1*, a nuclear receptor subfamily 4, group A gene (also known as *steroid receptor TR3* or *NUR77*) encodes an orphan member of the steroid-thyroid hormone-retinoid receptor superfamily whose members mainly act as transcriptional factors to positively or negatively regulate gene expression and play roles in regulating growth and apoptosis [28,29]. A role for *NR4A1* in cell proliferation has been previously reported. It was shown that its expression is rapidly induced by various mitogenic stimuli such as: serum growth factor, epidermal growth factor and fibroblast growth factor [28].

Taken together, the genes that we identified and validated appear to be involved in signaling pathways related to cell proliferation, invasion and metastatic processes, but their exact role in breast carcinogenesis remains to be elucidated.

#### Gene Ontology analysis

Classification of genes based on Gene Ontology (GO) terms is a powerful bioinformatics tool suited for the analysis of DNA microarray and SAGE data. Analysis of GO annotation allows one to identify families of genes that may play significant roles related to specific molecular or biological processes in expression profiles [30]. We used the *Expression Analysis Systematic Explorer* software (*EASE*) [31] to annotate the 520 deregulated genes according to the information provided by the GO Consortium [30]. The GO database provided annotation for 80% (419 out of 520) of the genes identified by SAGE. Results of this analysis are shown in Figure 2 and in detail in additional data file 2.

We observed that 31% of ER $\alpha$  associated transcripts are involved in biological processes related to *cell growth and/or maintenance*, 21% are related to *cell communication*, and 16% are related to *regulation of transcription*. Approximately 16% of these deregulated genes are related to molecular functions associated with *DNA binding* and more specifically with *transcription factor activity* (10%) (Figure 2). Interestingly, using the enrichment GO terms analysis, we identified statistical significant over-represen-



**Figure 2**  
**GO classification of the ER $\alpha$  associated genes identified by SAGE.** Percent of coverage representing the percentage of genes annotated with a specific GO term related to Biological Processes (blue bars) and Molecular Function (yellow bars).

tation of specific groups of proteins including: *metal ion binding proteins* (54 hits out of 419 annotated genes;  $p = 0.011$ ), *calcium ion binding proteins* (27 hits out of 419;  $p = 0.032$ ) and *steroid hormone receptor activity related proteins* (6 hits out of 419;  $p = 0.031$ ) (additional data file 2). The GO cluster related to *steroid hormone receptor activity proteins* includes: *estrogen receptor 1 (ESR1, i.e. ER $\alpha$ )*, *androgen receptor (AR)*, *hydroxysteroid 17- $\beta$  dehydrogenase 4 (HSD17 $\beta$ 4)*, *glucocorticoid receptor (NR3C1)*, *oxysterol binding protein (OSBP)*, and *retinoic acid receptor  $\alpha$  (RARA)*. The observation of functionally related groups of genes identified in the SAGE dataset via GO over representation analysis allows the identification of distinct biological pathways directly or indirectly associated to estrogen response related processes and provides the basis for future mechanistic studies.

#### Identification of high-affinity Estrogen Response Elements

We used a recently reported genome-wide high-affinity ERE database [32] to identify putative EREs in the promoter regions of the SAGE-identified 473 up-modulated genes in ER $\alpha$  (+) breast tumors. We identified 220 EREs distributed on 163 out of the 473 genes (35%) (see additional data file 3). Seventy-two percent of these genes contain one high affinity ERE (117 out of 163) and 28% of

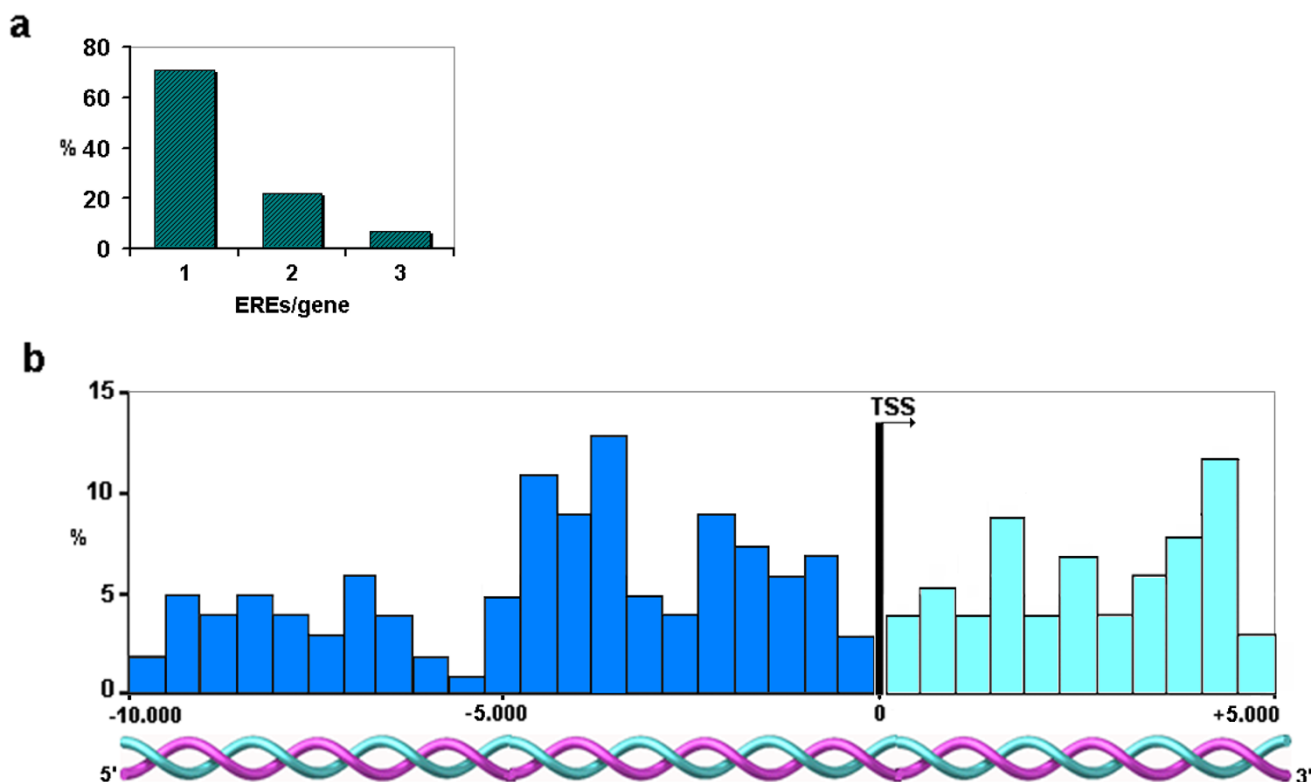
them contain two or more EREs in proximity to the transcriptional start sites (TSS) (46 out of 163) (Figure 3a). These EREs can be located in both coding and non-coding sequences such as was described by Bourdeau *et al.* [32].

The observed frequency of these elements in our study was 220 EREs in 3260 kb (considering a DNA window of 20 kb for each one of the 163 up-modulated genes with EREs). Compared with the expected frequency from random distribution of high-affinity EREs found in the genome (732 EREs in 3,069,334 kb 0.8 ERE in 3260 kb) (see material and methods) [32], the number of individual EREs was 270 fold higher than expected by chance ( $p < 0.00001$ ).

Fifty percent (110 out of 220) of the detected EREs mapped within a 10 kb region 5' of the TSS, while the rest mapped to 3' regions (Figure 3b). Approximately 68% of EREs mapped within the region between -5 to +5 kb from the TSS; in agreement with those observations of Bourdeau *et al.* [32]. However, it remains to be determined whether distantly located EREs (e.g. -10 kb from the TSS) are functional E<sub>2</sub>-ER binding sites related to transcriptional activation.

Of the validated transcripts previously discussed (Figure 1), we detected high-affinity EREs on the upstream or downstream regions related to the TSS of *SYTL4* (-8384 bp from the TSS: tggacatcatgacct), *TSPAN1* (+974 bp and +9384 bp from the TSS: tggctctgaatgacct and aggtcattccacct respectively), *CELSR2* (+173 bp and +3607 bp from the TSS: tgctcagggtgacct and aggtcaccatgaccg respectively), and *NR4A1* (-3478 bp and +4217 bp from the TSS: tgttcactctgacct).

It is interesting to note that we were unable to identify high-affinity EREs on the majority of deregulated genes (65%) associated with a positive ER $\alpha$  status. The possibility exists that many of these genes are transcriptionally regulated by non-ERE mediated mechanisms such as those involving ER binding to the AP1 or SP1 transcription factors [33]. The AP1 transcription factor is a heterodimer formed by Jun and Fos family member proteins that binds to the phorbol diester (TPA) response element as well as to the AP1 consensus DNA sequence. In this pathway, ER plays a co-activator role for AP1 [6]. The ER/AP-1 complex can confer estrogen responsiveness to additional subset of genes found in our dataset such as: *ovalbumin* (Fold change: 3;  $p = 0.033$ ) and *c-fos* (Fold change: 2.1;  $p = 0.033$ ); two transcripts detected as over-expressed in ER $\alpha$  (+) breast tumors by SAGE (additional data file 1). Similarly the ER/SP1 complex confers estrogen responsiveness to genes such as: *retinoic acid receptor  $\alpha$  (RARA)* (Fold change: 6.7;  $p = 0.038$ ), *vascular endothelial growth factor (VEGFC)* (Fold change: 2.6;  $p = 0.037$ ), *insu-*



**Figure 3**  
**High-affinity EREs in ER $\alpha$  (+) up-modulated genes (n = 163).** **a)** Percentage of genes according to number of EREs. **b)** Distribution of EREs in 5' (blue bars) and 3' (aquamarine bars) regions relative to the TSS (-10 to +5 kb). Each bar represents an interval width of 500 bp.

*lin-like growth factor binding protein-4 (IGFBP4)* (Fold change: 2;  $p = 0.01$ ) and *heat shock protein 27 (HSPB1)* (Fold change: 2;  $p = 0.045$ ); four transcripts detected as over-expressed in ER $\alpha$  (+) tumors in our study (additional data file 1).

An additional pathway of transcription regulation by estrogen involves the ER-related receptors (ERR), *nuclear orphan receptors* with significant homology to ERs, which do not bind estrogen and have unknown physiological ligands. ERRs are known to bind to the steroidogenic factor 1 response element (SFRE) and also bind to classic EREs, by means of which they exert constitutive transcriptional activity [34]. We detected over-expression of the nuclear orphan receptor *NR4A1* by SAGE and subsequently validated this observation by real time RT-PCR (Figure 1g). Interestingly, and as previously mentioned, the genomic region 5' and 3' to the TSS of *NR4A1* contain high-affinity EREs. Interaction between ERs and ERRs has been observed in the transcriptional regulation of certain genes such as the human breast cancer related gene *TFF1/*

*pS2*, the promoter of which is not only activated by ERs but also by ERRs [35].

As described, ER $\alpha$  can mediate estrogenic response through multiple genomic and non-genomic mechanisms, many of which affect proteins and pathways not necessarily directly or exclusively associated with ER $\alpha$ . Thus it is worth stressing that it will be the totality of deregulated proteins the ones that ultimately define the phenotype of ER $\alpha$  (+) breast carcinomas regardless of whether a "direct association" with ER transcriptional regulation exists or not.

#### **In vivo versus in vitro estrogen induced global gene expression findings**

The SAGE profiles for E<sub>2</sub>-responsive genes in MCF-7 cell line, previously reported by us [36], was compared with the ER status genes expression profile found in primary breast carcinomas. Briefly, we detected 199 transcripts differentially expressed ( $p < 0.01$ ) in MCF-7 treated cells, 124 were up-regulated and 75 were down-regulated

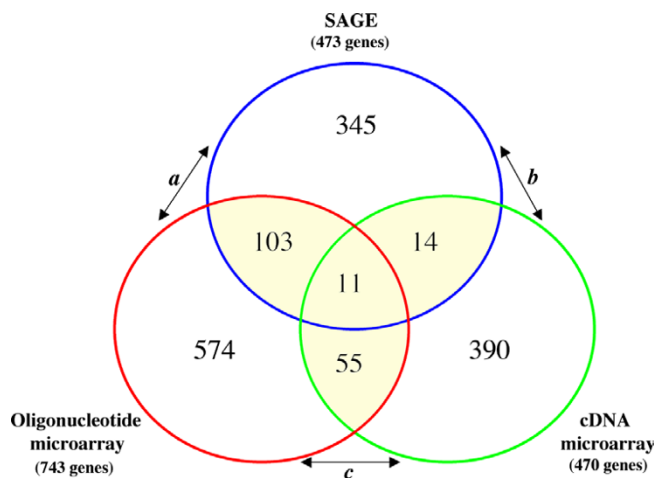
transcripts. Basically and as reported Charpentier *et al*, we observed a general up-regulation cell cycle progression-related genes including: *CCT2*, *CCND1*, *PES1*, *RAN/TC4*, *CALM1*, *CALM2*; and tumor-associated genes such as: *RFP*, *D52L1*, *TFF1/PS2*, *CAV1*, and *NDKA* among others [36]. These together could contribute to the stimulation of proliferation and the suppression of apoptosis by E<sub>2</sub>-ER transcriptional regulation.

By comparing the *in vitro* (199 differentially expressed transcripts) and *in vivo* (520 differentially expressed transcripts) gene expression profiles, to our surprise we detect that only few transcripts: *TFF1*, *CCND1*, *H19*, *SREBF1* and *WWP1* behaved similarly (i.e. up-regulation) in both studies. This is similar to observations made previously by Meltzer and co-workers whom showed that the majority of genes regulated in cell culture do not predict ER status in breast carcinomas [11,37]. This result suggests that the estrogen-responsive pathways affected *in vitro* represent only a minor portion of the global gene expression profiles characteristic of ER $\alpha$  (+) breast tumors. This maybe in great part the result of the heterogenous nature of bulk tumor tissue but in addition, the *in vitro* response of a single cell line to E<sub>2</sub>, in this particular case the widely used MCF-7 cells, may not faithfully reproduce the physiological effects of ER signaling *in vivo*.

#### Cross-platform gene expression profiling comparison

In order to identify and validate the most reliable set of genes able to discriminate breast carcinomas based on their ER $\alpha$  status, we performed a cross-platform comparison between the described SAGE dataset with two previously reported breast cancer studies based on DNA microarray methods [12,13]. van't Veer *et al*. [12] reported the gene expression profile of 97 primary breast tumors based on oligonucleotide microarrays containing 24,479 elements (Agilent Technologies, Palo Alto, CA, USA). In another study, Sotiriou *et al*. [13] reported the gene expression profile of 99 primary breast tumors using a cDNA microarray containing 7650 elements. Only files containing differentially expressed genes associated to ER $\alpha$  status tumors from both microarrays studies were obtained for cross-platform comparison (see material and methods).

Among the three platforms, a total of 1686 transcripts were identified as over-expressed in ER $\alpha$  (+) breast tumors. One hundred and eighty-three genes were identified by more than one method (Figure 4; additional data file 4). Eleven of these 183 genes were identified by all three methods displaying over-expression in ER $\alpha$  (+) breast carcinomas: *estrogen receptor 1 (ESR1)*, *GATA-binding protein 3 (GATA3)*, *mucin 1 (MUC1)*, *v-myb-myeloblastosis viral oncogene homolog (C-MYB)*, *X-box-binding protein 1 (XBP1)*, *hydroxysteroid 17- $\beta$  dehydrogenase 4 (HSD17B4)*,



**Figure 4**  
**Cross-platform comparisons of the up-modulated transcripts in ER $\alpha$  (+) breast carcinomas.** One hundred and eighty-three genes were identified by more than one study, eleven of which were commonly identified across the three platforms. **a)** Comparison between SAGE and oligonucleotide microarray platforms [12] showing a highly significant number of overlapping genes ( $p < 0.001$ ) (see table 2). **b)** Comparison between SAGE and cDNA microarray platforms [13] ( $p > 0.05$ ). **c)** Statistically significant number of overlapping genes identified by both DNA microarrays platforms ( $p < 0.01$ ).

*BTG family member 2 (BTG2)*, *transforming growth factor  $\beta$  3 (TGFB3)*, *member RAS oncogene family (RAB31)*, *START domain containing 10 (STARD10)*, and *KIAA0089* (Table 2).

One hundred and fourteen genes were identified as over-expressed by oligonucleotide microarrays [12] and SAGE in ER $\alpha$  (+) tumors, representing a non-random significant number of overlapping genes based on normal approximation to the binomial distribution ( $p < 0.001$ ) (Figure 4 and Table 2). Sixty-six genes were identified as over-expressed in ER $\alpha$  (+) tumors by both DNA microarrays platforms ( $p < 0.01$ ). The set of 25 genes overlapping between cDNA microarrays [13] and SAGE were not statistically significant ( $p > 0.05$ ).

Interestingly, we found a higher number of overlapping genes between the oligonucleotide microarray and SAGE platforms (114 genes), while only 66 genes were observed overlapping when comparing both microarray platforms. It is worth noting that 96% of the 470 genes (Figure 4) identified as overexpressed by the cDNA microarray method [13] were included within the total set of elements in the oligonucleotide microarray platform [12]. In



**Table 2: Transcripts identified as over-expressed in ER $\alpha$  (+) breast cancers commonly detected by cross-platforms comparison (SAGE and oligonucleotide microarrays).**

Gene name	Locus Link ID	Fold change	Frequency	Gene name	Locus Link	Fold change	Frequency#
<i>TFF1</i> *	7031	51.4	15/18 (83%)	<i>SULF2</i>	55959	2.9	11/18 (61%)
<i>SYTL4</i> *	94121	28.0	15/18 (83%)	<i>THBS4</i>	7060	2.9	8/18 (44%)
<i>DUSP4</i>	1846	14.7	14/18 (78%)	<i>AZGP1</i>	563	2.8	9/18 (50%)
<i>NAT1</i>	9	11.7	15/18 (83%)	<i>BBC3</i> *	27113	2.8	12/18 (67%)
<i>ECM1</i> *	1893	10.1	13/18 (72%)	<i>NET7</i> *	23555	2.8	10/18 (55%)
<i>KIAA0882</i>	23158	10.1	18/18 (100%)	<i>NET6</i>	27075	2.8	12/18 (67%)
<i>JDPI</i>	56521	10.0	12/18 (67%)	<i>TRAF5</i>	7188	2.8	9/18 (50%)
<i>ESRL</i>	2099	9.8	18/18 (100%)	<i>BTG2</i>	7832	2.7	9/18 (50%)
<i>HDGFRP3</i>	50810	9.8	12/18 (67%)	<i>RNF123</i> *	63891	2.7	11/18 (61%)
<i>TCEAL1</i>	9338	9.8	13/18 (72%)	<i>CHAD</i> *	1101	2.6	12/18 (67%)
<i>TSPAN1</i> *	10103	9.5	15/18 (83%)	<i>CSNK1A1</i>	1452	2.6	14/18 (78%)
<i>C20orf103</i> *	24141	9.3	10/18 (55%)	<i>EVL</i>	51466	2.6	12/18 (67%)
<i>MYLIP</i>	29116	9.3	11/18 (61%)	<i>HIST1H2BD</i>	3017	2.6	10/18 (55%)
<i>ABCA3</i>	21	8.9	10/18 (55%)	<i>SUSD3</i>	203328	2.6	9/18 (50%)
<i>SEC14L2</i>	23541	8.7	9/18 (50%)	<i>PLAT</i> *	5327	2.6	8/18 (44%)
<i>ANXA9</i> *	8416	8.4	10/18 (55%)	<i>RARRES3</i> *	5920	2.6	11/18 (61%)
<i>KCTD3</i>	51133	8.4	17/18 (94%)	<i>SH3BGR1</i> *	6451	2.6	8/18 (44%)
<i>SCUBE2</i>	57758	7.5	14/18 (78%)	<i>TPBG</i> *	7162	2.6	9/18 (50%)
<i>ITGBL1</i>	9358	7.1	8/18 (44%)	<i>UGCG</i>	7357	2.6	11/18 (61%)
<i>C14orf168</i>	83544	6.7	6/18 (33%)	<i>CELSR2</i> *	1952	2.5	8/18 (44%)
<i>FBP1</i>	2203	6.7	14/18 (78%)	<i>CRIMI</i>	51232	2.5	11/18 (61%)
<i>MYB</i>	4602	6.7	14/18 (78%)	<i>FLJ90798</i> *	219654	2.5	9/18 (50%)
<i>RARA</i> *	5914	6.7	12/18 (67%)	<i>KIF12</i>	113220	2.5	7/18 (39%)
<i>CaMKIIN<math>\alpha</math></i>	55450	6.3	18/18 (100%)	<i>LRIG1</i>	26018	2.5	9/18 (50%)
<i>AR</i> *	367	6.2	10/18 (55%)	<i>LRP2</i> *	4036	2.5	10/18 (55%)
<i>ZNF552</i>	79818	6.2	16/18 (89%)	<i>PHF15</i> *	23338	2.5	12/18 (67%)
<i>MIPEP</i> *	4285	6.0	14/18 (78%)	<i>HSMNP1</i>	55861	2.4	8/18 (44%)
<i>BAI2</i>	576	5.3	15/18 (83%)	<i>LOC123169</i>	123169	2.4	12/18 (67%)
<i>DP1L1</i>	92840	5.3	15/18 (83%)	<i>PINK1</i> *	65018	2.4	11/18 (61%)
<i>VAV3</i>	10451	5.3	12/18 (67%)	<i>PRKAR2B</i>	5577	2.4	7/18 (39%)
<i>KIAA0089</i>	23171	5.2	17/18 (94%)	<i>TJP3</i> *	27134	2.4	11/18 (61%)
<i>GATA3</i>	2625	5.1	15/18 (83%)	<i>CCND1</i>	595	2.3	9/18 (50%)
<i>QDPR</i>	5860	5.1	11/18 (61%)	<i>CYBRD1</i>	79901	2.3	10/18 (55%)
<i>C1orf21</i>	81563	4.9	11/18 (61%)	<i>KRT18</i>	3875	2.3	10/18 (55%)
<i>KIAA1143</i>	57456	4.9	7/18 (39%)	<i>PURA</i>	5813	2.3	9/18 (50%)
<i>OIP106</i>	22906	4.9	16/18 (89%)	<i>SREBF1</i> *	6720	2.3	10/18 (55%)
<i>AGR2</i>	10551	4.6	10/18 (55%)	<i>CYB5R1</i>	51706	2.2	6/18 (33%)
<i>MGC4251</i>	84336	4.6	13/18 (72%)	<i>DLG3</i> *	1741	2.2	9/18 (50%)
<i>FER1L3</i>	26509	4.4	10/18 (55%)	<i>EEF1A2</i>	1917	2.2	11/18 (61%)
<i>C4A</i>	720	4.1	11/18 (61%)	<i>GSTZ1</i>	2954	2.2	9/18 (50%)
<i>CRIP2</i>	1397	4.0	15/18 (83%)	<i>LOC159090</i>	159090	2.2	6/18 (33%)
<i>NTN4</i>	59277	4.0	10/18 (55%)	<i>MGC11242</i> *	79170	2.2	10/18 (55%)
<i>GJAI</i>	2697	3.8	11/18 (61%)	<i>MGC18216</i> *	145815	2.2	8/18 (44%)
<i>CGI-111</i> *	51015	3.7	14/18 (78%)	<i>NEIL1</i>	79661	2.2	6/18 (33%)
<i>CROT</i> *	54677	3.6	15/18 (83%)	<i>XBPI</i> *	7494	2.2	8/18 (44%)
<i>DACH</i>	1602	3.6	13/18 (72%)	<i>IRX5</i>	10265	2.1	8/18 (44%)
<i>DKFZP564D172</i>	83989	3.6	10/18 (55%)	<i>RAB31</i>	11031	2.1	9/18 (50%)
<i>FGD3</i>	89846	3.6	10/18 (55%)	<i>SSBP2</i>	23635	2.1	7/18 (39%)
<i>RNASE4</i> *	6038	3.6	12/18 (67%)	<i>TGFB3</i>	7043	2.1	8/18 (44%)
<i>GLUL</i> *	2752	3.3	11/18 (61%)	<i>BMPRI1B</i>	658	2.0	7/18 (39%)
<i>FOXAI</i>	3169	3.2	10/18 (55%)	<i>FLJ21174</i>	79921	2.0	6/18 (33%)
<i>MGC7036</i>	196383	3.2	14/18 (78%)	<i>FLJ22386</i>	79641	2.0	7/18 (39%)
<i>MUC1</i> *	4582	3.2	12/18 (67%)	<i>HSPB1</i> *	3315	2.0	6/18 (33%)
<i>NAVI</i>	89796	3.1	13/18 (72%)	<i>IGFBP4</i> *	3487	2.0	8/18 (44%)

**Table 2: Transcripts identified as over-expressed in ER $\alpha$  (+) breast cancers commonly detected by cross-platforms comparison (SAGE and oligonucleotide microarrays).** (Continued)

<b>RPLP1*</b>	6176	3.1	12/18 (67%)	<b>MGC15737*</b>	85012	2.0	8/18 (44%)
<b>ALCAM</b>	214	2.9	9/18 (50%)	<b>SPARCLI</b>	8404	2.0	9/18 (50%)
<b><u>HSD17B4*</u></b>	3295	2.9	13/18 (72%)	<b><u>STARD10*</u></b>	10809	2.0	7/18 (39%)

\* Genes with known or putative high-affinity EREs mapping in the vicinity of the TSS.

# Transcripts tags changing > 2-fold when compared with the average expression of ER $\alpha$  (-) tumors. Underlined genes correspond to the transcripts cross-validated among all three compared platforms.

other words, it appears that a better correlation was observed between SAGE and oligonucleotide arrays, than between both DNA microarray methods.

## Conclusion

In summary, our comprehensive comparison of overlapping genes across different gene expression platforms provides validation for a significant number of transcripts identified as highly expressed in ER $\alpha$  (+) breast tumors. More importantly this analysis identifies the most promising biomarkers for further evaluation as ER $\alpha$  associated genes in breast cancer. Furthermore, the identified proteins may be of value as breast cancer prognostic indicators analyzed either as a group or individually. It is also likely that groups of co-regulated genes in ER $\alpha$  (+) breast cancers may be associated to the hormonal control of mammary epithelial cells growth and differentiation. Finally, a better understanding of the signaling networks controlled or associated with the estrogen response may lead to the identification of novel breast cancer therapeutic targets.

## Methods

### SAGE libraries

To perform the comparative breast cancer SAGE analysis based on ER $\alpha$  status, we analyzed 26 Stage I – Stage II invasive breast carcinomas (8 ER $\alpha$ -negative tumors and 18 ER $\alpha$ -positive tumors). To this end, we generated and sequenced 24 breast cancer SAGE libraries at an approximate resolution of 100,000 tags per library, combined with 2 additional breast cancer libraries (ER $\alpha$ -negative tumors) downloaded from the Cancer Genome Anatomy Project – SAGE Genie database (SAGE\_Breast\_Carcinoma\_B\_95-259 and B\_IDC\_4) <http://cgap.nci.nih.gov/SAGE/>. For the generation of our SAGE libraries, snap frozen samples were obtained from the M.D. Anderson breast cancer tumor bank, and SAGE analysis was performed as previously described [36,38].

### Data processing and statistical analysis of SAGE libraries

SAGE tag extraction from sequencing files was performed by using the SAGE2000 software version 4.0 (a kind gift of Dr. K. Kinzler, John Hopkins University). SAGE data management, tag to gene matching as well as additional gene

annotations and links to publicly available resources such as GO, UniGene, LocusLink, were performed using a suite of web-based SAGE library tools developed by us [http://spi.mdacc.tmc.edu/bitools/about/sage\\_lib\\_tool.html](http://spi.mdacc.tmc.edu/bitools/about/sage_lib_tool.html). In our analyses we only considered tags with single tag-to-gene reliable matches. To compare these SAGE libraries, we utilized a modified t-test recently developed by us [18]. This test is based on a beta binomial sampling model that takes into account both, the intra-library and the inter-library variability, thus identifying 'common patterns' of SAGE transcript tag changes systematically occurring across samples [18].

All raw SAGE data reported as Supplementary tables in this manuscript is publicly available at [http://science.park.mdanderson.org/ggeg/sage\\_Proj\\_9.htm](http://science.park.mdanderson.org/ggeg/sage_Proj_9.htm).

### Real Time RT-PCR analysis

Template cDNAs were synthesized on mRNAs isolated from an independent set of 36 Stage I – Stage II human breast carcinomas (13 ER $\alpha$ -negative tumors and 23 ER $\alpha$ -positive tumors) obtained from our tumor bank. Primers and probes were obtained from the TaqMan Assays-on-Demand™ Gene Expression Products (Applied Biosystems, Foster City, CA, USA). All the PCR reactions were performed using the TaqMan PCR Core Reagents kit and the ABI Prism® 7700 Sequence Detection System (Applied Biosystems, Foster City, CA, USA). Experiments were performed in duplicate for each data point and 18s rRNA was used as control. Results were expressed as mean  $\pm$  2 Standard Error based on Log<sub>2</sub> transformation of normalized real time RT-PCR values of the assayed genes. We used t-test to compare the gene expression levels of validated genes between ER $\alpha$  (+) and ER $\alpha$  (-) breast tumors ( $p < 0.05$ ).

### Immunohistochemical determination of ER status

IHC staining and ER status determination was performed by the Pathology Department, MDACC following routine immunohistochemical procedures. Briefly, five micrometer sections of invasive breast carcinomas paraffin embedded tissues were used. Endogenous peroxidase activity was blocked with 3% H<sub>2</sub>O<sub>2</sub> in methanol for 10 min. After pretreatment with Tris-EDTA buffer, in order to block

non-specific antibody binding, the slides were incubated with 10% goat serum in PBS for 30 min. Primary monoclonal ER $\alpha$  antibody (ER-6F11, Novocastra, Newcastle, UK) was used at 1:50 dilution and detected following standard immunohistochemical techniques. DAB was used as chromogen and Mayers hematoxylin is used as counterstain. Scoring was performed by breast pathologist (AS). Cutoff for positivity was determined at 5% of tumor cells staining positively for ER (i.e. < 5% of cells the tumor was considered negative for ER $\alpha$ ).

#### Bioinformatics analysis

For automated functional annotation and classification of genes of interest based on GO terms, we used the EASE [31] available at the *Database for Annotation, Visualization and Integrated Discovery (DAVID)* at <http://david.niaid.nih.gov/david> [39]. The EASE software calculates over-representation of specific GO terms with respect to the total number of genes assayed and annotated. Statistical measures of specific enrichment of GO terms are determined by means of an EASE score ( $p < 0.05$ ). The EASE score is a conservative adjustment of the Fisher exact probability that weights significance in favor of biological themes supported by more genes and is calculated using the Gaussian hypergeometric probability distribution that describes sampling without replacement from a finite population [31]. This allows one to identify biological themes within a specific list of EASE analyzed genes.

#### High-affinity Estrogen Response Elements (ERE) analysis

To identify the occurrence of EREs within the promoter regions of up-modulated genes in ER $\alpha$  (+) breast tumors, we used a human genome-wide high-affinity ERE database <http://mapageweb.umontreal.ca/maders/eredatabase/> [32]. This public available database contains 71,119 EREs identified across the human genome (related to 17,353 transcriptional start sites), representing the consensus ERE (5'-Pu-GGTCA-NNN-TGACC-Py-3'), and equivalent sequences with only one or two nucleotide variations from such consensus. Based on these restrictions the expected random frequency was calculated as the total number of base pairs in the human genome divided by the frequency of occurrence of a sequence with specified base pairs at 10 positions and two base pair choices at two positions ( $3,069,334,246/4^{11} = 732$  high-affinity EREs) [32].

#### Comparison of gene expression patterns identified by different methodologies

ER $\alpha$  status associated genes identified in previous breast cancer studies [12,13] using DNA microarray methods were compared with our SAGE findings.

All over-expressed genes in ER $\alpha$  (+) breast tumors obtained from these studies were downloaded from the

corresponding web sites ([http://www.nature.com/cgi-taf/DynaPage.taf?file=/nature/journal/v415/n6871/abs/415530a\\_fs.html](http://www.nature.com/cgi-taf/DynaPage.taf?file=/nature/journal/v415/n6871/abs/415530a_fs.html) and <http://www.pnas.org/cgi/content/abstract/100/18/10393>) [12,13].

These datasets were annotated by LocusLink ID using the EASE software [26], and then compiled into one Excel spreadsheet pivotTable for comparison of overlapping genes between platforms, i.e. SAGE, Oligonucleotide and cDNA arrays. Anonymous ESTs from the microarrays platforms were excluded due to the inability to cross validate the identities between different gene expression profiles. Any combination of two lists was compared for matching gene-identity. The number and identity of genes commonly affected in two platforms (e.g. SAGE study vs. DNA microarray) was determined. We used the normal approximation to the binomial distribution as previously described [40] to calculate whether the number of matching genes derived from each cross-platform comparison was of statistical significance ( $p < 0.05$ ).

#### Authors' contributions

M.C.A. conceived the study idea and carried out the real time RT-PCR validations, the biostatistical/ bioinformatics analysis and writing the manuscript. Y.H., H.S. and J.A.D. carried out the breast cancer SAGE libraries and provided practical feedback on aspects of the manuscript. K.B. and S.G. developed the biostatistical and web-page base methodology. A.S. provides the tissue samples and clinical information. C.M.A. is the principal investigator and was involved in the conceptualization, design and writing of the manuscript. All authors read and approved the final manuscript.

#### Competing interests

The author(s) declare that they have no competing interests.

#### Additional material

##### Additional File 1

Differentially expressed genes between ER $\alpha$  (+) vs. ER  $\alpha$  (-) breast carcinomas (Fold change  $\geq 2$ ;  $p < 0.05$ ).

Click here for file

[<http://www.biomedcentral.com/content/supplementary/1471-2164-6-37-S1.xls>]

##### Additional File 2

Gene Ontology overrepresentation analysis.

Click here for file

[<http://www.biomedcentral.com/content/supplementary/1471-2164-6-37-S2.xls>]

**Additional File 3**

High-affinity EREs identified in ER $\alpha$  (+) up-modulated genes.

Click here for file

[<http://www.biomedcentral.com/content/supplementary/1471-2164-6-37-S3.xls>]

**Additional File 4**

Cross-platform comparison of the up-modulated transcripts in ER $\alpha$  (+) breast carcinomas.

Click here for file

[<http://www.biomedcentral.com/content/supplementary/1471-2164-6-37-S4.xls>]

**Acknowledgements**

The authors thank Dr. Michael MacLeod for critical reading of this manuscript. This work was supported by NIH-NCI Grant 1U19 CA84978-1A1 (C. M. Aldaz) and center grant ES-07784.

**References**

- Gruber CJ, Tschugguel W, Schneberger C, Huber JC: **Production and action of estrogens.** *N Engl J Med* 2002, **346**:340-352.
- Hall JM, Couse JF, Korach KS: **The multifaceted mechanisms of estradiol and estrogen receptor signaling.** *J Biol Chem* 2001, **276**:1642-1644.
- O'lone R, Frith MC, Karlsson EK, Hansen U: **Genomic targets of nuclear estrogen receptors.** *Mol Endocrinol* 2004, **18**:1859-1875.
- Klinge CM: **Estrogen receptor interaction with estrogen response elements.** *Nucleic Acids Res* 2001, **29**:2905-2919.
- Qin C, Singh P, Safe S: **Transcriptional activation of insulin-like growth factor-binding proteing-4 by 17 $\beta$ -estradiol in MCF-7 cell: role of estrogen receptor-Sp1 complexes.** *Endocrinology* 1999, **140**:2501-2508.
- Kushner PJ, Agard DA, Greene GL, Scanlan TS, Shiao AK, Uht RM, Webb P: **Estrogen receptor pathways to AP-1.** *J Steroid Biochem Mol Biol* 2000, **74**:311-317.
- Pedram A, Razandi M, Aitkenhead M, Hughes CCW, Levin ER: **Integration of the non-genomic and genomic actions of estrogen membrane-initiated signaling by steroid to transcription and cell biology.** *J Biol Chem* 2002, **277**:50768-50775.
- Shek LL, Dodolphin W: **Survival with breast cancer: the importance of estrogen receptor quantity.** *Eur J Cancer Clin Oncol* 1989, **25**:243-250.
- Pearce ST, Jordan VC: **The biological role of estrogen receptors  $\alpha$  and  $\beta$  in cancer.** *Crit Rev Oncol Hematol* 2004, **50**:3-22.
- Perou CM, Sorlie T, Eisen MB, van de Rijn M, Jeffrey SS, Rees CA, Pollack JR, Ross DT, Johnsen H, Akslen LA, Fluge O, Pergamenschikov A, Williams C, Zhu SX, Lonning PE, Borresen-Dale A, Brown PO, Botstein D: **Molecular portraits of human breast tumors.** *Nature* 2000, **406**:747-752.
- Gruvberger S, Ringner M, Chen Y, Panavally S, Saal LH, Borg A, Ferno M, Peterson C, Meltzer PS: **Estrogen receptor status in breast cancer is associated with remarkably distinct gene expression patterns.** *Cancer Res* 2001, **61**:5979-5984.
- van't Veer LJ, Dai H, van de Vijver MJ, He YD, Hart AAM, Mao M, Petersen HL, van der Kooy K, Marton MJ, Witteveen AT, Schreiber GJ, Kerkhoven RM, Roberts C, Linsley PS, Bernards R, Friend SH: **Gene expression profiling predicts clinical outcome of breast cancer.** *Nature* 2002, **415**:530-536.
- Sotiriou C, Neo S, McShane LM, Korn EL, Long PM, Jazaeri A, Martiat P, Fox SB, Harris AL, Liu ET: **Breast cancer classification and prognosis based on gene expression profiles from a population-based study.** *Proc Natl Acad Sci USA* 2003, **100**:10393-10398.
- Sorlie T, Tibshirani R, Parker J, Hastie T, Marron JS, Nobel A, Deng S, Johnsen H, Pesich R, Geisler S, Demeter J, Perou CM, Lonning PE, Brown PO, Borresen-Dale A, Botstein D: **Repeated observation of breast tumor subtypes in independent gene expression data sets.** *Proc Natl Acad Sci USA* 2003, **100**:8418-8423.
- West M, Blanchette C, Dressman H, Huang E, Ishida S, Spang R, Zuzan H, Olson JA, Marks JR, Nevins JR: **Predicting the clinical status of human breast cancer by using gene expression profiles.** *Proc Natl Acad Sci USA* 2001, **98**:11462-11467.
- Liu ET, Sotiriou C: **Defining the galaxy of gene expression in breast cancer.** *Breast Cancer Res* 2002, **4**:141-144.
- Velculescu VE, Zhang L, Vogelstein B, Kinzler KW: **Serial analysis of gene expression.** *Science* 1995, **270**:484-487.
- Baggerly KA, Deng L, Morris JS, Aldaz CM: **Differential expression in SAGE: accounting for normal between-library variation.** *Bioinformatics* 2003, **19**:1477-1483.
- Yang R, Domingos CK, Wasserman SM, Colman SD, Shenoy S, Mehraban F, Komuves LG, Tomlinson JE, Topper JN: **Identification of a novel family of cell-surface proteins expressed in human vascular endothelium.** *J Biol Chem* 2002, **277**:46364-46373.
- Grimmond S, Larder R, Van Hateren N, Siggers P, Hulsebos TJM, Arkell R, Greenfield A: **Cloning, mapping, and expression analysis of gene encoding a novel Mammalian EGF-related protein (SCUBE1).** *Genomics* 2000, **70**:74-81.
- Usui T, Shima Y, Shimada Y, Hirano S, Burgess RW, Schwarz TL, Takeichi M, Uemura T: **Flamingo, a seven-pass transmembrane cadherin, regulates planar cell polarity under the control of Fz/Dishevelled.** *Cell* 1999, **98**:585-595.
- Formstone CJ, Little PFR: **The flamingo-related mouse Celsr family (Celsr1-3) genes exhibit distinct patterns of expression during embryonic development.** *Mech Develop* 2001, **109**:91-94.
- Shima Y, Kengaku M, Hirano T, Takeichi M, Uemura T: **Regulation of dendritic maintenance and growth by a mammalian 7-pass transmembrane cadherin.** *Development* 2004, **131**:205-216.
- Kuroda TS, Fukuda M, Ariga H, Mikoshiba K: **The Slp homology domain of synaptotagmin-like proteins 1-4 and Slac2 functions as a novel Rab27A binding domain.** *J Biol Chem* 2002, **277**:9212-9218.
- Hao X, Sun B, Hu L, Lahdesmaki H, Dunmire V, Feng Y, Zhang SW, Wang H, Wu C, Wang H, Fuller GN, Symmans WF, Shmulevich I, Zhang W: **Differential gene and protein expression in primary breast malignancies and their lymph node metastases as revealed by combined cDNA microarray and tissue microarray analysis.** *Cancer* 2004, **100**:1110-1122.
- Berditchevski F, Odintsova E: **Characterization of integrin-tetraspanin adhesion complexes: role of tetraspanins in integrin signaling.** *J Cell Biol* 1999, **146**:477-492.
- Sugiura T, Berditchevski F: **Function of  $\alpha$ 3 $\beta$ 1-Tetraspanin protein complex in tumor cell invasion evidence for the role of the complexes in production of matrix metalloproteinase 2 (MMP-2).** *J Cell Biol* 1999, **146**:1375-1389.
- Kolluri SK, Bruey-Sedano N, Cao X, Lin B, Lin F, Han Y, Dawson MI, Zhang X: **Mitogenic effect of orphan receptor TR3 and its regulation by MEK1 in lung cancer.** *Mol Cell Biol* 2003, **23**:8651-8667.
- Lin B, Kolluri SK, Lin F, Liu W, Han Y, Cao X, Dawson MI, Reed JC, Zhang X: **Conversion of Bcl-2 from protector to killer by interaction with nuclear orphan receptor Nur77/TR3.** *Cell* 2004, **116**:527-540.
- Ashburner M, Ball CA, Blake JA, Botstein D, Butler H, Cherry JM, Davis AP, Dolinski K, Dwight SS, Eppig JT, Harris MA, Hill DP, Issel-Tarver L, Kasarskis A, Lewis S, Matese JC, Richardson JE, Ringwald M, Rubin GM, Sherlock G: **Gene ontology: tool for the unification of biology. The Gene Ontology Consortium.** *Nat Genet* 2000, **25**:25-29.
- Hosack DA, Dennis G, Sherman BT, Lane HC, Lempicki RA: **Identifying biological themes within lists of genes with EASE.** *Genome Biol* 2003, **4**:R70.
- Bourdeau V, Deschenes J, Metivier R, Nagai Y, Nguyen D, Bretschneider N, Gannon F, White JH, Mader S: **Genome-wide identification of high-affinity estrogen response elements in human and mouse.** *Mol Endocrinol* 2004, **18**:1411-1427.
- Gruber CJ, Gruber DM, Gruber IML, Wieser F, Huber JC: **Anatomy of the estrogen response element.** *Trends Endocrinol Metab* 2004, **15**:73-78.
- Vanacker JM, Bonnelye E, Chopin-Delannoy S, Delmarre C, Cavailles V, Laudet V: **Transcriptional activities of the orphan nuclear receptor EER alpha (estrogen receptor-related receptor-alpha).** *Mol Endocrinol* 1999, **13**:764-773.

35. Lu D, Kiriya Y, Lee KY, Giguere V: **Transcriptional regulation of estrogen-inducible pS2 breast cancer marker gene by the ERR family of orphan nuclear receptors.** *Cancer Res* 2001, **61**:6755-6761.
36. Charpentier AH, Bednarek AK, Daniel RL, Hawkins KA, Laflin KJ, Gaddis S, MacLeod MC, Aldaz CM: **Effects of estrogen on global gene expression: identification of novel targets of estrogen action.** *Cancer Res* 2000, **60**:5977-5983.
37. Cunliffe HE, Ringner M, Bilke S, Walker RL, Cheung JM, Chen Y, Meltzer PS: **The gene expression response of breast cancer to growth regulators: patterns and correlation with tumor expression profiles.** *Cancer Res* 2003, **63**:7158-7166.
38. Abba MC, Drake JA, Hawkins KA, Hu Y, Sun H, Notcovich C, Gaddis S, Sahin A, Baggerly K, Aldaz CM: **Transcriptomic changes in human breast cancer progression as determined by serial analysis of gene expression.** *Breast Cancer Res* 2004, **6**:R499-R513.
39. Dennis G, Sherman BT, Hosack DA, Yang J, Gao W, Lane HC, Lempicki RA: **DAVID: Database for Annotation, Visualization, and Integrated Discovery.** *Genome Biol* 2003, **4**:R60.
40. Smid M, Dorssers LCJ, Jenster G: **Venn Mapping: clustering of heterologous microarray data based on the number of co-occurring differentially expressed genes.** *Bioinformatics* 2003, **19**:2065-2071.

Publish with **BioMed Central** and every scientist can read your work free of charge

*"BioMed Central will be the most significant development for disseminating the results of biomedical research in our lifetime."*

Sir Paul Nurse, Cancer Research UK

Your research papers will be:

- available free of charge to the entire biomedical community
- peer reviewed and published immediately upon acceptance
- cited in PubMed and archived on PubMed Central
- yours — you keep the copyright

Submit your manuscript here:  
[http://www.biomedcentral.com/info/publishing\\_adv.asp](http://www.biomedcentral.com/info/publishing_adv.asp)

