# Gene expression profiles of poor-prognosis primary breast cancer correlate with survival

François Bertucci[1,2,7], Valéry Nasser[1], Samuel Granjeaud[6], François Eisinger[3], José Adelaïde[1], Rebecca Tagett[4], Béatrice Loriod[6], Aurélia Giaconia[1], Athmane Benziane[4], Elisabeth Devilard[5], Jocelyne Jacquemier[5], Patrice Viens[2,7], Catherine Nguyen[6], Daniel Birnbaum[1,7,8,*] and Rémi Houlgatte[6]

[1]Département d'Oncologie Moléculaire TAGC, [2]Département d'Oncologie Médicale, [3]Département de Prévention et Dépistage, [4]Ipsogen SA and [5]Département d'Anatomie Pathologique, Institut Paoli-Calmettes, [6]TAGC, CIML Luminy, [7]Université de la Méditerranée and [8]Laboratoire d'Oncologie Moléculaire, U.119 Inserm, Marseille, France

The extensive heterogeneity of breast cancer complicates the precise assessment of tumour aggressiveness, making therapeutic decisions difficult and treatments inappropriate in some cases. Consequently, the long-term metastasis-free survival rate of patients receiving adjuvant chemotherapy is only 60%. There is a genuine need to identify parameters that might accurately predict the effectiveness of this treatment for each patient. Using cDNA arrays, we profiled tumour samples from 55 women with poor-prognosis breast cancer treated with adjuvant anthracycline-based chemotherapy. Gene expression monitoring was applied to a set of about 1000 candidate cancer genes. Differences in expression profiles provided molecular evidence of the clinical heterogeneity of disease. First, we confirmed the capacity of a 23-gene predictor set, identified in a previous study, to distinguish between tumours associated with different survival. Second, using a refined gene set derived from the previous one, we distinguished, among the 55 clinically homogeneous tumours, three classes with significantly different clinical outcome: 5-year overall survival and metastasis-free survival rates were respectively 100% and 75% in the first class, 65% and 56% in the second and 40% and 20% in the third. This discrimination resulted from the differential expression of two clusters of genes encoding proteins with diverse functions, including the estrogen receptor (ER). Another finding was the identification of two ER-positive tumour subgroups with different survival. These results indicate that gene expression profiling can predict clinical outcome and lead to a more precise classification of breast tumours. Furthermore, the characterization of discriminator genes might accelerate the development of new specific and alternative therapies, allowing more rationally tailored treatments that are potentially more efficient and less toxic.

## INTRODUCTION

Over the last decades, advances in systemic adjuvant therapy, designed to eradicate the micrometastases observed at diagnosis, have substantially improved the treatment of poor-prognosis primary breast cancer (1,2). But conventional clinicopathological factors remain insufficient to evaluate the substantial prognostic heterogeneity of this disease. Therefore, a recent consensus conference led by the US National Cancer Institute (NCI) has recommended the enlargement of the criteria leading to the use of adjuvant chemotherapy, which still relies on clinical and morphological parameters and the use of anthracycline-based regimens as standard treatment (http://odp.od.nih.gov/consensus/). Consequently, more and more patients are offered this therapy [with a 60% long-term metastasis-free survival rate (1)] whose success or failure in an individual patient cannot currently be predicted by any clinical or morphological factor. Alternative therapeutic strategies have been developed, and have already had promising results in advanced disease. Among those being tested as adjuvant treatments are new cytotoxic agents (e.g. taxanes), new hormonal therapies and new biological agents (e.g. trastuzumab) (3). There is clearly a crucial need to identify parameters that might accurately predict the clinical outcome after specific adjuvant treatment in individual patients, allowing a rational choice between the different therapies available, improving efficiency and reducing morbidity and cost of treatment.

*To whom correspondence should be addressed at: Laboratoire d'Oncologie Moléculaire, U.119 Inserm, IFR57, 27 Boulevard Leï Roure, 13009 Marseille, France. Tel: +33 4 91 75 84 07; Fax: +33 4 91 26 03 64; Email: birnbaum@marseille.inserm.fr

Breast cancer is a multifactorial disease characterized by the accumulation of numerous molecular alterations in the cells. This complexity makes each tumour potentially distinct from all others at the molecular and clinical levels. Prognosis and resistance to treatment are not likely to be associated with the disturbance of a single gene, but rather with the combined influence of many genes. Comprehensive molecular analyses should help identify such genes, and cDNA array technology (4,5), which allows the analysis of the mRNA expression levels of thousands of genes simultaneously in a sample, could be the method of choice. Tumour cell models have suggested the utility of this approach for investigating major issues in metastasis (6,7) and chemoresistance (8–10). Notably, the response of cancer cell lines to certain drugs has been shown to be predictable by gene expression signatures (9). Several recent studies have demonstrated the usefulness of expression profiles for improving the classification of cancers by identifying new subgroups of tumours within clinically and morphologically similar groups (11–19). Such a molecular taxonomy has suggested prognostic information for lymphomas (12), renal cell carcinomas (20), and oesophageal (21), lung (22,23) and breast cancers (14,24).

Our aim is to identify, within apparently homogeneous populations of samples, new previously unrecognized tumour classes displaying distinct clinical courses after therapy. By analysing a limited series of primary breast carcinomas with cDNA arrays, we previously identified, using a supervised analysis, a predictor set of 23 genes whose expression patterns differentiated two groups of patients with different survival after adjuvant chemotherapy (14). To validate and further extend these results, we present here the expression analysis of some 1000 candidate genes from a larger, independent and homogeneous series of poor-prognosis primary breast cancers treated with adjuvant chemotherapy. We confirm the prognostic classification provided by the previously identified predictor set of 23 genes. Then we improve this predictor set and refine the tumour classification by sorting the samples into three classes with significantly different long-term survival.

## RESULTS

### Gene expression profiling of breast cancer

The mRNA from 66 different human breast cancer samples, including 55 clinical tissue samples and 11 cell lines, were hybridized with cDNA arrays containing about 1000 selected genes. The overall expression patterns for these 66 samples were analysed with hierarchical clustering and displayed in a colour-coded matrix (Fig. 1). The clustering algorithm classifies samples on the horizontal axis and genes on the vertical axis, ordered on the basis of similarity of their expression profiles. Overall similarity of breast tumours and cell lines is shown as a dendrogram where branch length reflects relatedness of the samples (Fig. 1A). This analysis highlighted groups of correlated genes across correlated samples (Fig. 1B). Some interesting gene clusters are indicated by coloured bars on the left of Fig. 1B and are shown enlarged in Fig. 1C. Three of these were differentially expressed between tissue samples and epithelial cell lines. A 'stromal cluster' (blue

bar) and an 'immune cluster' (green bar) were overexpressed in tissues overall as compared with cell lines, probably reflecting the inflammatory component of the tumours. These clusters were rich in genes whose expression is respectively found in stromal cells (collagen genes COL1A1, COL6A1, proteases MMP2, MMP3, microfibrillar-associated protein 2 MFAP2) and in B cells (immunoglobulin genes, CD79, HLA class II, . . .), T-cells (CD2, CD3, TRB, TRD, . . .), monocytes or macrophages (CD14, CSF1R). The third cluster ('proliferation cluster', pink bar) included several genes involved in cell proliferation such as CDK4, ODC1, GSTP1, UBCH10, DNMT1, TUBA1, HDAC2 and PCNA (which codes for a proliferation marker used in clinical practice). This cluster was overexpressed in cell lines overall as compared with tissues, probably reflecting the difference in proliferation rate between rapidly dividing cells in culture and asynchronously proliferating cells in tumour tissues.

Comparison between tumour tissue RNA indicated a great heterogeneity of expression profiles, which were distinct for each tumour. The samples were classified in three large branches by the clustering algorithm (Fig. 1A, C). The 'stromal', 'immune' and 'proliferation' clusters were differentially expressed in these three categories. Similar grouping was observed using other classification algorithms such as k-means. As expected, the 'immune cluster' and the 'proliferation cluster' were overexpressed in estrogen receptor (ER)-negative tumours overall as compared with ER-positive tumours, in agreement with respectively a more abundant inflammatory stroma and a higher proliferation index (25). Another differentially expressed gene cluster contained the immediate-early genes FOS, JUNB and EGR1, which code for transcription factors involved in signalling pathways triggered by proliferation stimuli ('immediate response cluster', red bar in Fig. 1B, C). Interestingly, this cluster has previously been observed in breast (13) and ovarian (18) cancers. Another cluster (designated cluster I, orange bar in Fig. 1B, C) included ESR1, which codes for the estrogen receptor (ER) several transcription factor genes (GATA3, ILF1, XBP1, CRABP2, SMARCA2, ELF1, BS69 and GLI3) and the anti-apoptotic gene BCL2. Variation in expression of ESR1 mRNA correlated well with that of the protein measured by immunohistochemistry (IHC; concordance in 50 out of 55 samples). Finally, a cluster (designated cluster II, brown bar in Fig. 1B, C) was highly overexpressed in the middle group of tumours. It contained genes coding for membrane proteins (connexin 43/GJA1, cadherin 15, prolactin receptor PRLR, endothelin receptor EDNRA, mucin-like hormone receptor EMR1, P-glycoprotein 1/ABCB1/MDR1, MLANA/MART1) and genes involved in the cell cycle or apoptosis (CDKN3, DAP3, CIDEA).

### Classification of breast cancer and selection of gene clusters with potential prognostic role

As shown in Fig. 1A, classification with the whole set of genes identified three large groups of tumours differing with respect to their IHC ER status, but similar with respect to histological type or clinical outcome. The absence of prognostic classification could be due to the presence of irrelevant clusters of co-expressed genes that exert a dominant influence upon the clustering.
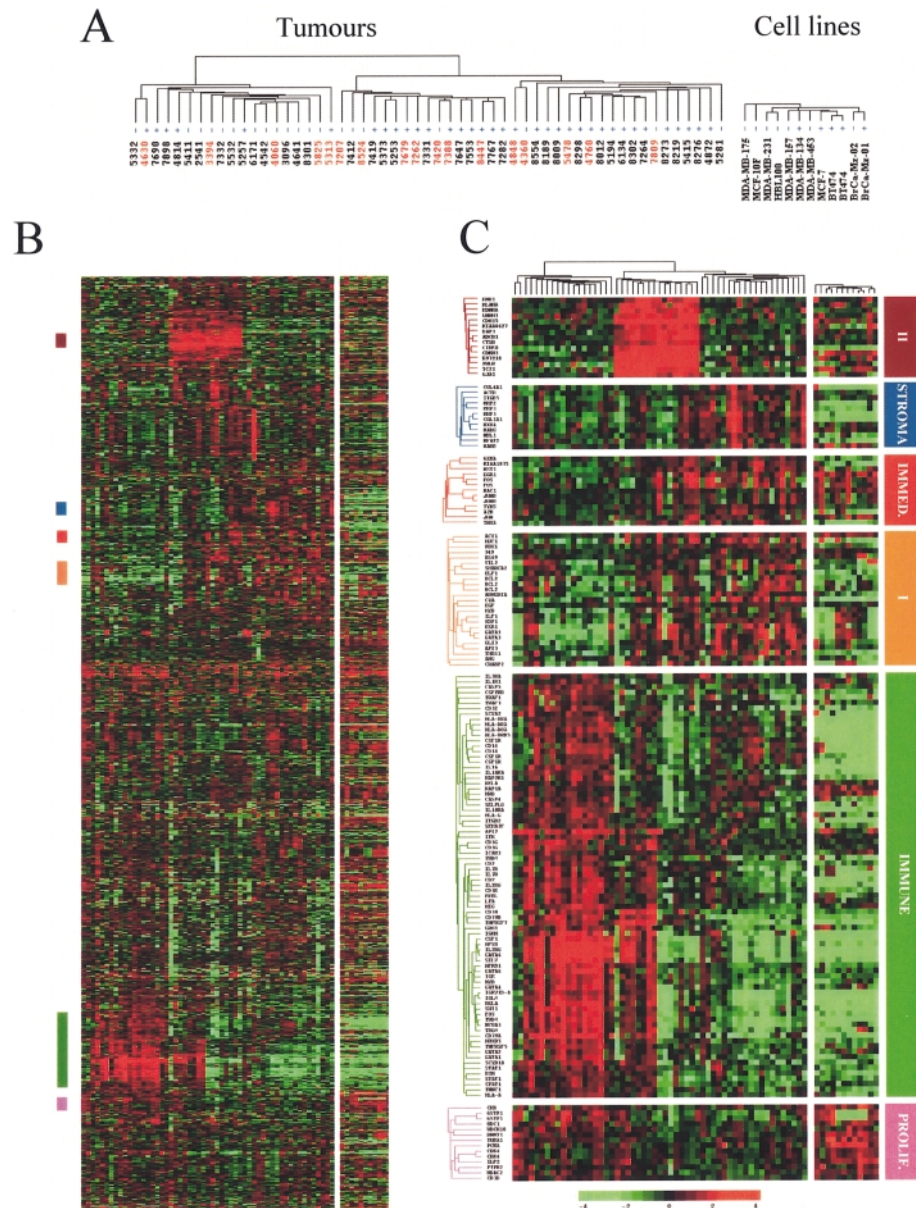
Figure 1. Expression patterns of 1045 cDNA clones in 66 experimental breast cancer samples. Each row represents a gene and each column represents a sample. Genes are referenced by their HUGO abbreviation as used in 'Locus Link' (http://www.ncbi.nlm.nih.gov/LocusLink/). Tumour tissue samples are designated with numbers. Each cell in the matrix represents the expression level of a transcript in a single sample relative to its median abundance across all samples and is depicted according to a colour scale shown at the bottom. Red and green indicate expression levels respectively above and below the median. The magnitude of deviation from the median is represented by the colour saturation. A hierarchical clustering was applied to group genes on the basis of similarity of their expression patterns across all samples. The same clustering was then separately applied to cell lines and tissue samples to group them on the basis of the similarity of their expression patterns. (A) Dendrogram of samples representing overall similarities in gene expression profiles across all samples. ER status measured by IHC is indicated for each sample ( + , positive;  − , negative). Fatal tumours are coloured red. (B) Matrix representation of expression levels. Coloured bars to the left indicate the locations of gene clusters of interest shown in (C), which is an expanded view of selected gene clusters named from top to bottom: cluster II, 'stromal cluster' (stroma), 'immediate response cluster' (immed.), cluster I, 'immune cluster' (immune) and 'proliferation cluster' (prolif.).

We first measured the predictive power of the expression levels of a set of 23 discriminator genes, which we previously identified using a supervised analysis in a study of 34 different tumours (14). The clustering identified two groups of patients (Fig. 2A) with significantly divergent 5-year survival rates: 53% in the left group and 87% in the right group (P < 0.05, log-rank test). The composition of the gene clusters differen-tially expressed in the good-prognosis group and in the poor-prognosis group was globally similar to that of our previous study. These results confirmed and reinforced our previous study with a larger and independent series of tumours.

Because of some degree of residual clinical heterogeneity in the two groups of tumours, we wanted to refine our classification. The objective was to identify, starting from the
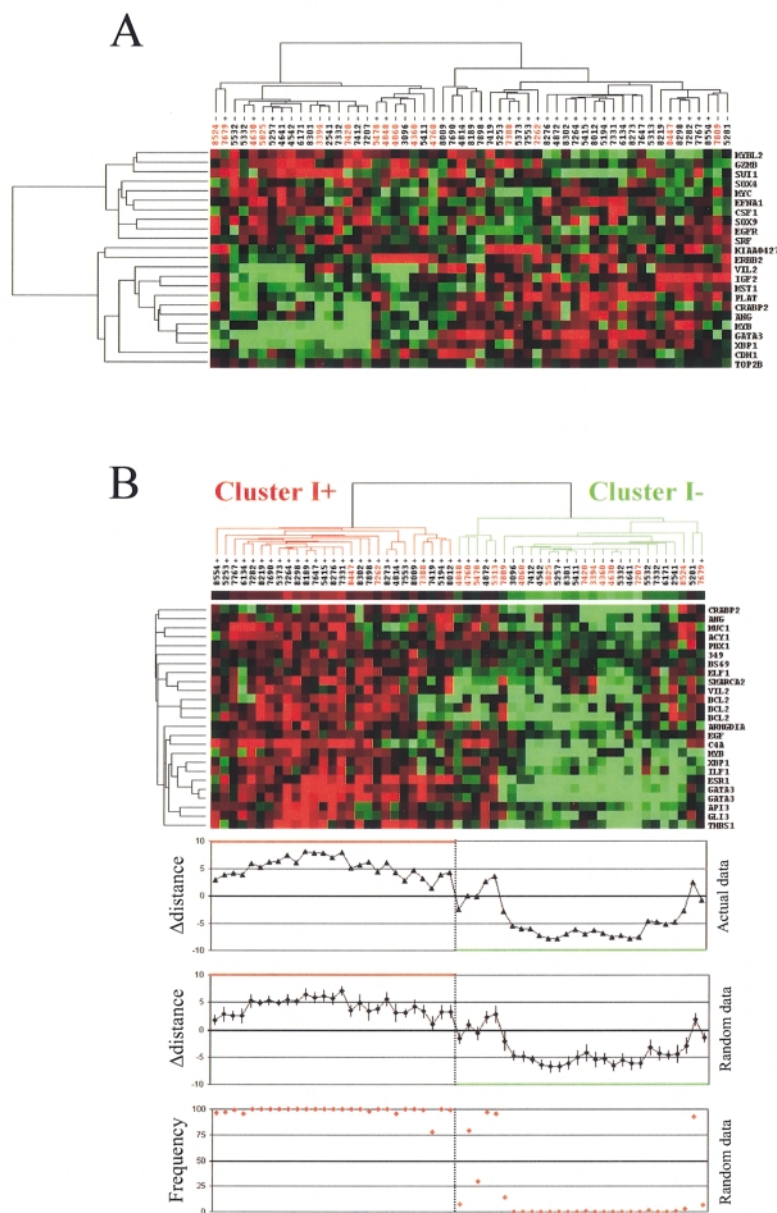
Figure 2. Classification of breast tumours using the 23-gene set and cluster I. The 55 breast cancer samples were clustered using expression levels of two subsets of genes. (A) Hierarchical clustering based on the 23 previously identified discriminator genes (14). (B) Reclustering based on the cluster I genes identified in Fig. 1C. Two large groups of tumours were separated (red branches and green branches). Boxes under the coloured matrix display the accuracy of the clustering-based tumour classification, estimated by measuring the Euclidean distance D between each sample and the green (GD) and the red (RD) mean expression profiles (see Materials and Methods). The difference of Euclidean distance, $\Delta$distance ($\Delta D = GD - RD$), is plotted for each sample and reflects the proximity of its expression profile to each average profile. Samples are ordered according to their location in the original hierarchical clustering. The dotted vertical lines mark the separation between 'cluster I+' and 'cluster I−' tumour groups. Top box: $\Delta$distance estimated with actual data. Middle box: $\Delta$distance estimated with 10 000 randomly generated datasets (mean and standard deviation are respectively represented with black lozenges and error bars). Bottom box: frequency of membership of each sample to the red group measured with 10 000 randomly generated datasets.

23-gene cluster, particular gene subsets that would more precisely predict patient survival. The 23 genes did not form an individual cluster in the present study (data not shown): 9 were found in or proximal to cluster I, 4 were associated with the 'immune' and the 'proliferation' clusters, and 2 were associated with cluster II, while the remaining genes were scattered, suggesting an heterogeneous composition of this initial predictor set.

We first assessed whether cluster I alone (25 genes) was able to discriminate patients with different outcome. It contained ESR1, which encodes a known prognostic factor in breast cancer. Clustering all tissue samples using the expression levels of cluster I genes sorted two groups of tumours: 'cluster I+' and 'cluster I−' (respectively red and green branches in Fig. 2B). As expected, these groups differed in ESR1 mRNA expression, but the overlap was not perfect with IHC ER status.

While the 'cluster I+' group included only IHC-positive tumours, the 'cluster I−' group included both IHC-negative and IHC-positive tumours. The overall survival was significantly different between the two groups, with 3 deaths out of 27 patients in the 'cluster I+' group and 14 deaths out of 28 in the 'cluster I−' group (P < 0.05, log-rank test). Interestingly, patients with ER-positive 'cluster I−' tumours had a significantly shorter survival than patients with ER-positive 'cluster I+' tumours (P < 0.005, log-rank test). These results suggest that expression profiles of ESR1-associated genes provide different and more accurate clinical information than the IHC status alone, possibly reflecting functional differences in the ESR1 pathway.

Confidence for the two tumour groups revealed by clustering was assessed by measuring the distance between each sample and the mean profile of each group. As shown in Fig. 2B (top box), 50 samples were closer to their original group, while 3 were closer to the opposite group and 2 exhibited the same distance relative to the red or green groups. These 5 tumours displayed an 'intermediate cluster I profile'. Confidence was further checked by using randomly generated subsets of 25 genes from a larger set of genes correlated with ESR1 (using a correlation cut-off of 0.5 in the gene dendrogram). This resampling allowed the estimation of the mean and standard deviation of the differences of distances for each sample to the red and the green mean profiles (Fig. 2B, middle box) and the frequency of membership of each sample to the red group (Fig. 2B, bottom box). Results showed that 51 samples were initially well classified, while 4 'cluster I−' tumours were closer to the 'cluster I+' group. This classification discrepancy between clustering and distance to mean profile methods confirmed the 'intermediate cluster I profile' for these 4 tumours, perhaps due to a partially functional ESR1 pathway.

While these results confirmed the utility of our initial predictor gene set, as well as cluster I, to separate clinically identical tumours, the two groups of tumours still displayed some residual clinical heterogeneity (three deaths in the better-prognosis 'cluster I+' group).

### Gene-expression-based classification of breast cancer and survival

Focusing only on the 'cluster I+' tumours, we identified a gene cluster whose expression profile could separate two subgroups of samples, one of which included the three patients who died (Fig. 3A, B). Since this gene cluster was similar to cluster II, cluster II genes were used to recluster the 55 tumours. The two groups of tumours (Fig. 3C) that emerged had very distinct expression profiles (Fig. 3C, boxes) and different clinical outcomes (although not statistically significant: P = 0.11, log-rank test). Further clustering of all 55 samples with all genes from clusters I and II failed to identify groups with significantly different survival (data not shown).

A complementarity between clusters I and II was found, however, when the corresponding classifications were analysed in a '2D representation' (Fig. 4A). This representation delineated four groups of tumours, A, B, C and D, which were in close agreement with the clinical outcome of patients. Overexpression of cluster I together with underexpression of cluster II defined a group (A) with 'good prognosis' and no

mortality in 16 patients. Underexpression of cluster I together with overexpression of cluster II defined a group (D) with 'poor prognosis': 4 deaths out of 5 patients. The two other groups (B and C) defined 'intermediate prognosis', with respectively 3 deaths out of 11 women in B and 10 out of 23 in C. Groups B and C were merged (class B + C), thus defining three classes. With a median follow-up of 60 months, the difference for overall survival (OS) was statistically significant between the three classes (P < 0.005, log-rank test). Concerning metastases, the difference was also statistically significant (P < 0.05, log-rank test), with 4 relapses out of 16 in A, 15 out of 34 in B + C and 4 out of 5 in D. Figures 4B and C respectively show the Kaplan Meier plots of overall survival and metastasis-free survival (MFS) from these three classes of tumours. Although the 5-year OS and MFS rates were 72% and 58% respectively for the whole population, they were 100% and 75% in the first class (A), 65% and 56% in the second (B + C), and 40% and 20% in the last (D).

Interestingly, these three classes did not show any significant difference with respect to follow-up or to classical breast cancer prognostic factors such as axillary lymph node status, patient age and menopausal status, histological type, and grade and size of tumours. Furthermore, no such prognostic survival classification was possible by axillary lymph node (negative versus positive; P = 0.89 for OS and P = 0.49 for MFS, log-rank test) or IHC ER status (negative versus positive; P = 0.45 for OS and P = 0.88 for MFS, log-rank test), both of which are classical and strong prognostic factors of breast cancer. The same analysis restricted to the 42 ductal tumours (11 samples in A, 27 in B + C and 4 in D) gave similar results, with significantly different survivals between the three classes (data not shown). Finally, survival differences between the three classes remained significant (P < 0.05, log-rank test), even when all 5 of the 'intermediate cluster I profile' tumours were artificially relocated to another class (from C to A) or excluded from analysis, further indicating the robustness of our results.

### DISCUSSION

Several recent studies have suggested the usefulness of comprehensive cDNA array-based gene expression profiles for cancer classification (13,15–19), and some have directly addressed the issue of prognosis (12,14,20–24). Using this technology, we profiled a series of 55 primary breast cancers using a set of candidate genes, the large majority of which are directly or indirectly implicated in oncogenesis. All patients had a poor-prognosis tumour according to the criteria in use at our institute, and had received adjuvant anthracycline-based chemotherapy after surgery. Today, all would receive the same treatment according to the recent NCI consensus recommendations. Follow-up of patients was sufficiently long to consider survival (median, 5 years), which conformed to data in the literature for similar population and treatment (1).

The expression patterns generated delineated clinically relevant classes of tumours. We first confirmed our initial results with a 23-gene predictor set defined in a previously studied cohort of patients (14). This set identified, within the present larger and independent series of tumours, two groups with different survival rates. Then, the analysis of more genes
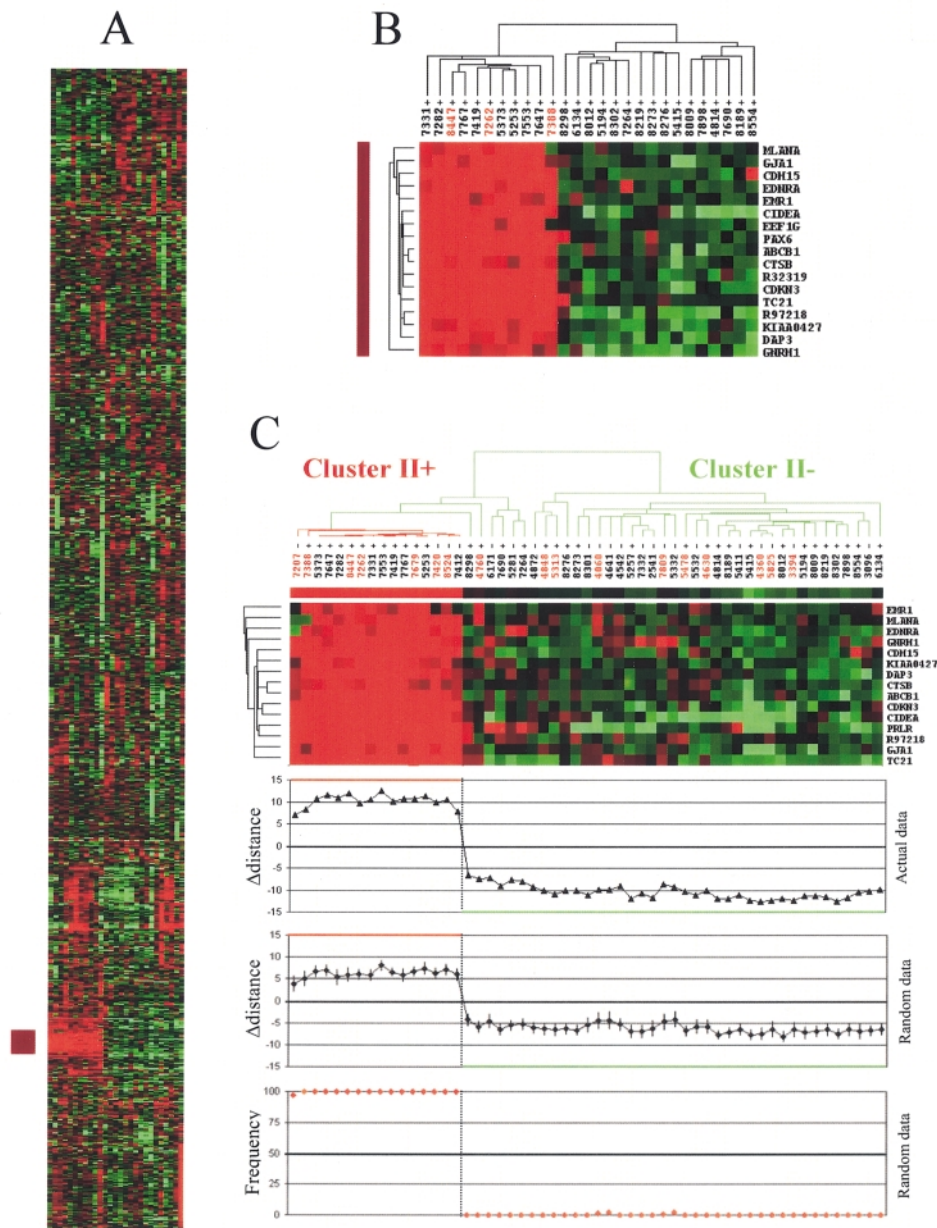
Figure 3. Cluster II and classification of breast tumour samples. (A), Hierarchical clustering of the 27 'cluster I+' tumours based on the whole set of genes. The brown bar to the left indicates the location of the gene cluster shown in (B). (B), The dendrogram at the bottom lists the 27 tumours and their classification as defined in (A). Note that the 3 patients who died are in the left group. The gene cluster responsible for this classification was compositionally very similar to cluster II. (C) Hierarchical clustering of 55 tumours based on the cluster II genes identified in Fig. 1C. Boxes under the coloured matrix display the accuracy of the clustering-based tumour classification (for legend, see Fig. 2C).

and samples allowed us to dissect this predictor gene cluster and to identify two derived gene clusters whose expression further refined the prognostic classification. One of the clusters, cluster I, included ESR1, a proven prognostic factor in breast cancer (26). With respect to clinical outcome, the expression signature of this cluster differentiated tumours better than the IHC ER status. It also revealed the existence of at least two molecular subtypes of ER-positive breast tumours (as measured by IHC) with different expression profiles and survival (longer survival in 'cluster I+' tumours). Subtypes of ER-positive breast tumours have recently been reported by others using

discriminator genes such as ESR1, GATA3, XBP1 or MYB, also present in our cluster I (24). The second gene cluster was only useful as a complement to cluster I. Their combined expression profiles drastically improved tumour classification and allowed the differentiation of three distinct classes with significantly different long-term survivals (P < 0.005 for OS and P < 0.05 for MFS). Interestingly, these classes were equilibrated with respect to clinicopathological features of samples. No such accurate classification could have been obtained using classical breast cancer prognostic parameters, suggesting that only gene expression profiles can distinguish relevant classes in this
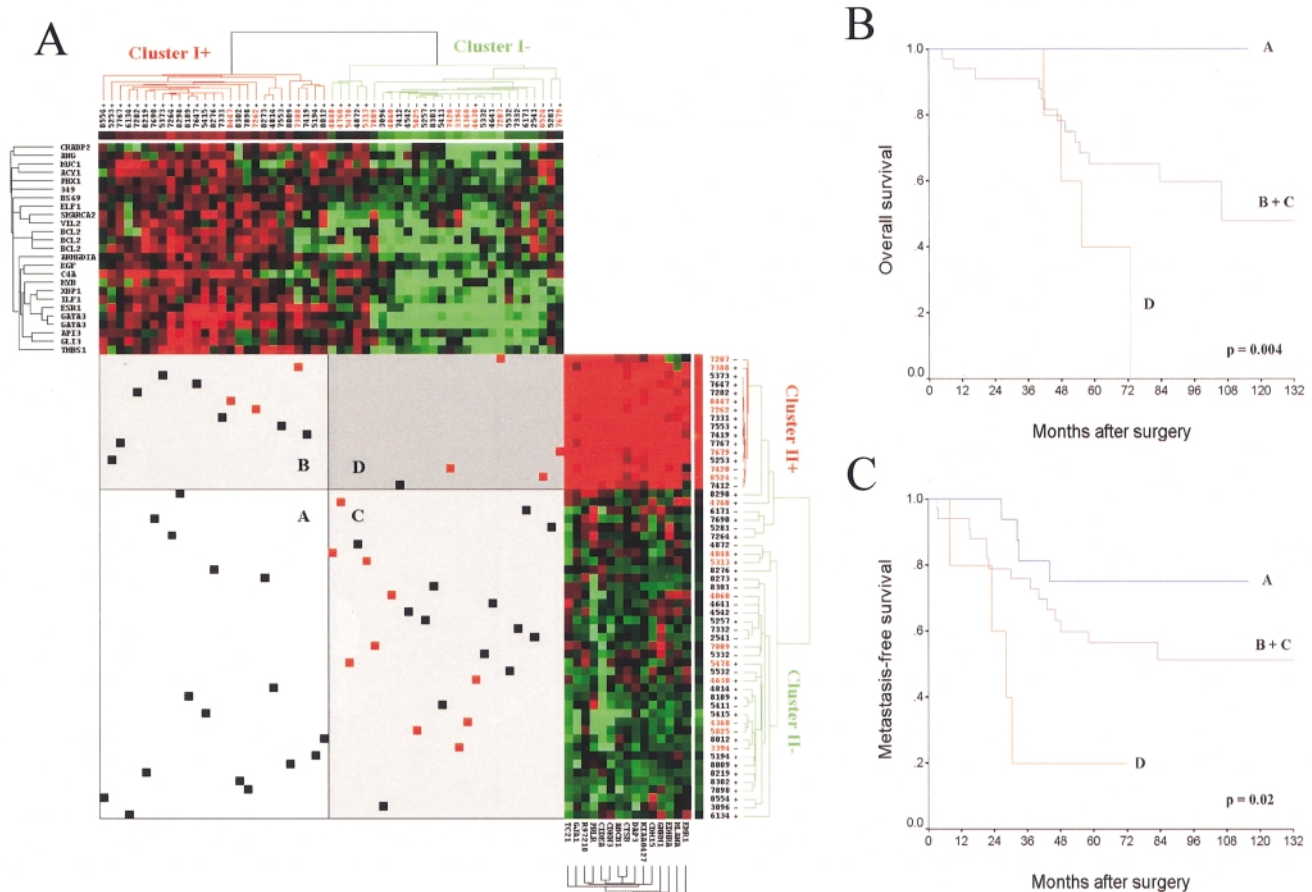
Figure 4. Gene-expression-based tumour classification correlates with clinical outcome. (A) Two-dimensional representation of the hierarchical clustering results are shown in Figs 2B and 3C. The analysis delineates four groups of tumours: A, B, C and D. Black squares indicate patients alive at last follow-up visit and red squares indicate patients who died. Three classes of patients with a statistically different clinical outcome were defined according to gene expression profiles: class A (n = 16), class B + C (n = 34), class D (n = 5). (B) Kaplan–Meier plot of overall survival of the three classes of patients (P < 0.005, log-rank test). (C) Kaplan–Meier plot of metastasis-free survival of the three classes of patients (P < 0.05, log-rank test).

apparently homogeneous prognostic category of patients, whose survival is in fact very different after chemotherapy. These results should have important therapeutic implications, guiding consensual poor-prognosis patients towards the type of adjuvant therapy most likely to succeed for them. Patients with tumours similar to class A would be candidates for standard adjuvant treatment, whereas an alternative therapy would have to be considered for the other patients. We must highlight however, that molecular heterogeneity – and to a lesser extent, clinical heterogeneity – is still present in each class, suggesting that more subtypes of breast tumours probably exist. Five tumours with the 'intermediate cluster I profile' probably represent another subtype of breast cancer, which might be distinguishable by profiling more samples.

The functional identities of the discriminator genes provide insight into mammary oncogenesis and may help identify potential therapeutic targets. Cluster I was overexpressed in the group of tumours with the best prognosis. It included ESR1 and several co-regulated genes. ER, a transcription factor, plays a critical role in the mammary gland, where it regulates cell proliferation, differentiation and motility in concert with other signalling pathways. Several genes of cluster I are known to be associated with ER in that they are estrogen-regulated and/or associated with a positive ER status. Examples are mucin 1 (MUC1) (27), the proto-oncogene MYB (28), BCL2 (29) and the cellular retinoic acid-binding protein 2 (CRABP2) (30). Transcriptional activators of ER are also included: the epidermal growth factor gene (EGF) (31), SMARCA2 (32), and the RHO–GDP dissociation inhibitor $\alpha$ (ARHGDIA) (6). Other transcription factors in cluster I include GATA3 (13,14,33), XBP1 (13,14) (whose mRNA expression has been previously associated with ER status), ILF1, ELF1, BS69, GLI3 and PBX1. The search for downstream genes regulated by these factors could help identify other genes potentially relevant to disease outcome. It is probable that the expression of this gene cluster reflects the functional status of the ER pathway, clinically more significant than the IHC ER status alone. This might explain the better survival observed for ER-positive tumours that overexpressed this gene cluster compared with ER-positive tumours which underexpressed it. Finally, this cluster included two genes involved in angiogenesis: angiogenin (ANG) (34) and thrombospondin 1 (THBS1) (35).

Cluster II was overexpressed in the group of tumours with the worst prognosis. It contained highly co-regulated genes with diverse functions such as oncogenes, hormone receptors and apoptotic factors. A notable feature was the inclusion of the ABCB1 gene coding for the membrane-associated P-glycoprotein 1, an ATP-binding cassette chloride channel involved in multidrug resistance and poor response to chemotherapy in breast cancer (36). Among the other genes were the cathepsin B protease gene (CTSB) implicated in metastasis (37), the RAS-like GTPase TC21 (38) and CDKN3, a cell cycle regulator (39). Interestingly, several hormone and hormone receptor transcripts were represented (PRLR, EMR1 and GNRH1). Prolactin (PRL) functions as an autocrine/paracrine factor that stimulates growth and differentiation of mammary tissue after binding to its receptor (PRLR). Although the activation of the PRL/PRLR pathway increases cellular motility in vitro (40), the role of PRL and PRLR in breast cancer progression remains poorly delineated (41). GJA1, the predominant connexin in human mammary epithelium, has been found to be overexpressed in doxorubicin-resistant breast cancer cell lines (10). Whichever might be the driving genes of these predictive clusters, their prognostic accuracy and value are certainly increased by the presence of other co-regulated genes, highlighting the utility of large-scale molecular analyses for explaining tumour heterogeneity. The identification of the most relevant genes in each class is interesting because they represent novel potential markers of sensitivity to current anticancer drugs and/or tumour aggressiveness. In addition, new specific therapies targeting these genes might be developed as alternatives to standard chemotherapy in the patients with poor outcome (classes B + C and D).

In conclusion, our study shows that gene expression profiles, defined using candidate gene arrays, can identify clinically relevant tumour subgroups, significantly contributing to the refinement of poor-prognosis breast cancer classification. By delineating discriminator genes, new alternative anticancer drugs might soon be developed. The application of better tailored and more specific therapy should lead to major improvements in cancer management, with better chances of success and potentially fewer side-effects. The next important step will be the analysis of larger series of patients in prospective clinical trials to assess the true impact of cDNA array data on patient treatment.

## MATERIALS AND METHODS

### Breast tumour samples and characteristics of patients

Tumour samples were obtained from 55 women treated at the Institut Paoli-Calmettes. These were chosen after careful screening based on the following criteria: (i) sporadic primary breast cancer treated with surgery followed by adjuvant anthracycline-based chemotherapy, (ii) tumour material quickly dissected and frozen in liquid nitrogen and stored at $-160^{\circ}$C, (iii) patient follow-up 48 months or more after diagnosis. In addition to the axillary node status, four poor-prognosis criteria were used to determine whether adjuvant chemotherapy should be administered: patient age less than 40 years, pathological tumour size greater than 20 mm, Scarff–Bloom–Richardson

(SBR) grade equal to 3, and negative ER status as evaluated by IHC (with a positivity cut-off value of 1%). Women who received chemotherapy were those either with node-positive tumours, or with node-negative tumours and one of the poor-prognosis criteria if non-menopausal or two criteria if menopausal. After surgery, all patients received comparable regimens of chemotherapy containing conventional doses of anthracycline every 21 days for six cycles. All tumour sections were de novo reviewed by a pathologist (J.J.) prior to analysis. All samples contained more than 50% tumour cells. Tumours were infiltrating adenocarcinomas, including (according to the WHO histological typing) 42 ductal carcinomas, 5 lobular, 5 mixed, and 3 medullary. Other main characteristics of patients are listed in Table 1.

### Breast cancer cell lines

All breast-cancer-derived cell lines were obtained from the American Type Culture Collection – except BrCa-MZ-01 and BrCA-MZ-02, which were kind gifts from Dr V.J. Möbus (Ulm, Germany) – and were grown as recommended.

### Complementary DNA array production

The 1045 human cDNA clones used in the study were obtained from the IMAGE consortium. These were chosen to represent genes with proven or suspected roles in cancer, including genes involved in transcription, the cell cycle, cell adhesion, invasion, angiogenesis and chemoresistance (the complete list is available at http:/tagc.univ-mrs.fr/pub/Cancer/). The use of control clones, PCR amplification and robotical spotting of PCR products onto Hybond-N+ membranes (Amersham) were done as described previously (42).

Table 1. Clinical characteristics of patients

| Characteristics | All patients (n = 55) |
|---|---|
| Median age (years) | 56 |
| Menopausal status (n): | |
| Postmenopausal | 30 |
| Premenopausal | 25 |
| Axillary lymph-node metastasis (n)[a]: | |
| Negative | 11 |
| Positive | 43 |
| Pathological tumour size (pT) (n): | |
| pT1 | 20 |
| pT2 | 27 |
| pT3 | 8 |
| SBR grade (n)[b]: | |
| I and II | 29 |
| III | 24 |
| Estrogen-receptor status (n): | |
| Positive | 35 |
| Negative | 20 |
| Median follow-up (months) | 60 |
| Metastatic relapse (n) | 23 |
| Death (n) | 17 |

[a]Not available in 1 patient.
[b]Not available in 2 patients.

### RNA extraction, hybridizations and data acquisition

Total RNA was extracted from frozen tumour samples and cell lines by standard methods (43). RNA integrity was controlled by 28S northern blots before labelling. Hybridizations of cDNA arrays were done with radioactive [$\alpha$-$^{33}$P]dCTP-labelled probes made from 5 µg of total RNA from each sample according to described protocols (http://tagc.univ-mrs.fr/pub/Cancer/). After washes, arrays were exposed to phosphoimaging plates, which were then scanned with a FUJI BAS 1500 machine.

### Data analysis and statistical methods

Signal intensities were quantified, normalized for the amount of spotted DNA (44) and the variability of experimental conditions, and log-transformed (42). Average-linkage hierarchical clustering was then applied to investigate relationships between samples and relationships between genes. We used the Cluster program (with Pearson correlation as similarity metric) and displayed results with the TreeView program (45).

The accuracy of tumour classification defined by hierarchical clustering was assessed. We determined the average expression profile of each tumour group defined by the original classification (green and red groups in Figs 2B and 3C) and measured the Euclidean distance D between each sample and the green (GD) and the red (RD) mean expression profiles. The difference between these distances ($\Delta$distance, $\Delta D = GD - RD$) reflected the proximity to average profiles ($\Delta D$ was positive for samples closer to the red profile and negative for samples closer to the green profile). Then, we used the intrinsic noise in the experimental data to estimate confidence levels for $\Delta D$ and the probability of each tumour belonging to each group. For every gene cluster tested, we retained a larger set of genes that contained, in addition to the original genes, their neighbour genes belonging to dendrogram branches with a correlation cut-off of 0.5. 10 000 datasets of the same size as the initial selected cluster were randomly generated from this new set, allowing, for each sample, the computation of $\Delta D$ (mean and standard deviation) and the frequency with which it was found closer to the red mean profile.

Survival analysis was done with the SPSS software (version 10.0.5). The primary endpoint was the overall survival (OS) measured from the time of diagnosis until the date of the last follow-up visit or cancer-related death. Metastasis-free survival (MFS) was the secondary endpoint and was measured in the same way until the date of the first distant metastasis. Survivals were estimated with the Kaplan–Meier method and compared between groups using the log-rank test (46). Data concerning patients who were alive or without metastatic relapse at last followup were censored. A P-value of less than 0.05 was considered significant.

## REFERENCES

1. Early Breast Cancer Trialists' Collaborative Group (1998) Polychemotherapy for early breast cancer: an overview of the randomised trials. Lancet, 352, 930–942.
2. Early Breast Cancer Trialists' Collaborative Group (1998) Tamoxifen for early breast cancer: an overview of the randomised trials. Lancet, 351, 1451–1467.
3. McCarthy, N.J. and Swain, S.M. (2000) Update on adjuvant chemotherapy for early breast cancer. Oncology (Huntingt), 14, 1267–1280; discussion 1280–1264, 1287–1268.
4. The Chipping Forecast (1999) Nat. Genet., 21(Suppl.), 1–60.
5. Granjeaud, S., Bertucci, F. and Jordan, B.R. (1999) Expression profiling: DNA arrays in many guises. Bioessays, 21, 781–790.
6. Clark, E.A., Golub, T.R., Lander, E.S. and Hynes, R.O. (2000) Genomic analysis of metastasis reveals an essential role for RhoC. Nature, 406, 532–535.
7. Zajchowski, D.A., Bartholdi, M.F., Gong, Y., Webster, L., Liu, H.L., Munishkin, A., Beauheim, C., Harvey, S., Ethier, S.P. and Johnson, P.H. (2001) Identification of gene expression profiles that predict the aggressive behavior of breast cancer cells. Cancer Res., 61, 5168–5178.
8. Scherf, U., Ross, D.T., Waltham, M., Smith, L.H., Lee, J.K., Tanabe, L., Kohn, K.W., Reinhold, W.C., Myers, T.G., Andrews, D.T. et al. (2000) A gene expression database for the molecular pharmacology of cancer. Nat. Genet., 24, 236–244.
9. Staunton, J.E., Slonim, D.K., Coller, H.A., Tamayo, P., Angelo, M.J., Park, J., Scherf, U., Lee, J.K., Reinhold, W.O., Weinstein, J.N. et al. (2001) Chemosensitivity prediction by transcriptional profiling. Proc. Natl Acad. Sci. USA, 98, 10787–10792.
10. Turton, N.J., Judah, D.J., Riley, J., Davies, R., Lipson, D., Styles, J.A., Smith, A.G. and Gant, T.W. (2001) Gene expression and amplification in breast carcinoma cells with intrinsic and acquired doxorubicin resistance. Oncogene, 20, 1300–1306.
11. Golub, T.R., Slonim, D.K., Tamayo, P., Huard, C., Gaasenbeek, M., Mesirov, J.P., Coller, H., Loh, M.L., Downing, J.R., Caligiuri, M.A. et al. (1999) Molecular classification of cancer: class discovery and class prediction by gene expression monitoring. Science, 286, 531–537.
12. Alizadeh, A.A., Eisen, M.B., Davis, R.E., Ma, C., Lossos, I.S., Rosenwald, A., Boldrick, J.C., Sabet, H., Tran, T., Yu, X. et al. (2000) Distinct types of diffuse large B-cell lymphoma identified by gene expression profiling. Nature, 403, 503–511.
13. Perou, C.M., Sorlie, T., Eisen, M.B., van de Rijn, M., Jeffrey, S.S., Rees, C.A., Pollack, J.R., Ross, D.T., Johnsen, H., Akslen, L.A. et al. (2000) Molecular portraits of human breast tumours. Nature, 406, 747–752.
14. Bertucci, F., Houlgatte, R., Benziane, A., Granjeaud, S., Adelaide, J., Tagett, R., Loriod, B., Jacquemier, J., Viens, P., Jordan, B. et al. (2000) Gene expression profiling of primary breast carcinomas using arrays of candidate genes. Hum. Mol. Genet., 9, 2981–2991.
15. Martin, K.J., Kritzman, B.M., Price, L.M., Koh, B., Kwan, C.P., Zhang, X., Mackay, A., O'Hare, M.J., Kaelin, C.M., Mutter, G.L. et al. (2000) Linking gene expression patterns to therapeutic groups in breast cancer. Cancer Res., 60, 2232–2238.
16. Hedenfalk, I., Duggan, D., Chen, Y., Radmacher, M., Bittner, M., Simon, R., Meltzer, P., Gusterson, B., Esteller, M., Kallioniemi, O.P. et al. (2001) Gene-expression profiles in hereditary breast cancer. N. Engl. J. Med., 344, 539–548.
17. Bittner, M., Meltzer, P., Chen, Y., Jiang, Y., Seftor, E., Hendrix, M., Radmacher, M., Simon, R., Yakhini, Z., Ben-Dor, A. et al. (2000) Molecular classification of cutaneous malignant melanoma by gene expression profiling. Nature, 406, 536–540.
18. Welsh, J.B., Zarrinkar, P.P., Sapinoso, L.M., Kern, S.G., Behling, C.A., Monk, B.J., Lockhart, D.J., Burger, R.A. and Hampton, G.M. (2001) Analysis of gene expression profiles in normal and neoplastic ovarian tissue

samples identifies candidate molecular markers of epithelial ovarian cancer. Proc. Natl Acad. Sci. USA, 98, 1176–1181.

19. Khan, J., Wei, J.S., Ringner, M., Saal, L.H., Ladanyi, M., Westermann, F., Berthold, F., Schwab, M., Antonescu, C.R., Peterson, C. et al. (2001) Classification and diagnostic prediction of cancers using gene expression profiling and artificial neural networks. Nat. Med., 7, 673–679.

20. Takahashi, M., Rhodes, D.R., Furge, K.A., Kanayama, H., Kagawa, S., Haab, B.B. and Teh, B.T. (2001) Gene expression profiling of clear cell renal cell carcinoma: gene identification and prognostic classification. Proc. Natl Acad. Sci. USA, 98, 9754–9759.

21. Kihara, C., Tsunoda, T., Tanaka, T., Yamana, H., Furukawa, Y., Ono, K., Kitahara, O., Zembutsu, H., Yanagawa, R., Hirata, K. et al. (2001) Prediction of sensitivity of esophageal tumors to adjuvant chemotherapy by cDNA microarray analysis of gene-expression profiles. Cancer Res., 61, 6474–6479.

22. Garber, M.E., Troyanskaya, O.G., Schluens, K., Petersen, S., Thaesler, Z., Pacyna-Gengelbach, M., van De Rijn, M., Rosen, G.D., Perou, C.M., Whyte, R.I. et al. (2001) Diversity of gene expression in adenocarcinoma of the lung. Proc. Natl Acad. Sci. USA, 13, 13784–13789.

23. Bhattacharjee, A., Richards, W.G., Staunton, J., Li, C., Monti, S., Vasa, P., Ladd, C., Beheshti, J., Bueno, R., Gillette, M. et al. (2001) Classification of human lung carcinomas by mRNA expression profiling reveals distinct adenocarcinoma subclasses. Proc. Natl Acad. Sci. USA, 13, 13790–13795.

24. Sorlie, T., Perou, C.M., Tibshirani, R., Aas, T., Geisler, S., Johnsen, H., Hastie, T., Eisen, M.B., van de Rijn, M., Jeffrey, S.S. et al. (2001) Gene expression patterns of breast carcinomas distinguish tumor subclasses with clinical implications. Proc. Natl Acad. Sci. USA, 98, 10869–10874.

25. Fisher, E.R., Osborne, C.K., McGuire, W.L., Redmond, C., Knight, W.A., 3rd, Fisher, B., Bannayan, G., Walder, A., Gregory, E.J., Jacobsen, A. et al. (1981) Correlation of primary breast cancer histopathology and estrogen receptor content. Breast Cancer Res. Treat., 1, 37–41.

26. Chang, J., Powles, T.J., Allred, D.C., Ashley, S.E., Clark, G.M., Makris, A., Assersohn, L., Gregory, R.K., Osborne, C.K., and Dowsett, M. (1999) Biologic markers as predictors of clinical outcome from systemic therapy for primary operable breast cancer. J. Clin. Oncol., 17, 3058–3063.

27. McGuckin, M.A., Walsh, M.D., Hohn, B.G., Ward, B.G., and Wright, R.G. (1995) Prognostic significance of MUC1 epithelial mucin expression in breast cancer. Hum. Pathol., 26, 432–439.

28. Guerin, M., Sheng, Z.M., Andrieu, N., and Riou, G. (1990) Strong association between c-myb and oestrogen-receptor expression in human breast cancer. Oncogene, 5, 131–135.

29. Binder, C., Marx, D., Overhoff, R., Binder, L., Schauer, A. and Hiddemann, W. (1995) Bcl-2 protein expression in breast cancer in relation to established prognostic factors and other clinicopathological variables. Ann. Oncol., 6, 1005–1010.

30. Bucco, R.A., Zheng, W.L., Wardlaw, S.A., Davis, J.T., Sierra-Rivera, E., Osteen, K.G., Melner, M.H., Kakkad, B.P. and Ong, D.E. (1996) Regulation and localization of cellular retinol-binding protein, retinol-binding protein, cellular retinoic acid-binding protein (CRABP) and CRABP II in the uterus of the pseudopregnant rat. Endocrinology, 137, 3111–3122.

31. Kato, S., Endoh, H., Masuhiro, Y., Kitamoto, T., Uchiyama, S., Sasaki, H., Masushige, S., Gotoh, Y., Nishida, E., Kawashima, H. et al. (1995) Activation of the estrogen receptor through phosphorylation by mitogen-activated protein kinase. Science, 270, 1491–1494.

32. Chiba, H., Muramatsu, M., Nomoto, A. and Kato, H. (1994) Two human homologues of Saccharomyces cerevisiae SWI2/SNF2 and Drosophila brahma are transcriptional coactivators cooperating with the estrogen receptor and the retinoic acid receptor. Nucleic Acids Res., 22, 1815–1820.

33. Hoch, R.V., Thompson, D.A., Baker, R.J. and Weigel, R.J. (1999) GATA-3 is expressed in association with estrogen receptor in breast cancer. Int. J. Cancer, 84, 122–128.

34. Montero, S., Guzman, C., Cortes-Funes, H. and Colomer, R. (1998) Angiogenin expression and prognosis in primary breast carcinoma. Clin. Cancer Res., 4, 2161–2168.

35. Weinstat-Saslow, D.L., Zabrenetzky, V.S., VanHoutte, K., Frazier, W.A., Roberts, D.D. and Steeg, P.S. (1994) Transfection of thrombospondin 1 complementary DNA into a human breast carcinoma cell line reduces primary tumor growth, metastatic potential, and angiogenesis. Cancer Res., 54, 6504–6511.

36. Trock, B.J., Leonessa, F. and Clarke, R. (1997) Multidrug resistance in breast cancer: a meta-analysis of MDR1/gp170 expression and its possible functional significance. J. Natl Cancer Inst., 89, 917–931.

37. Foekens, J.A., Kos, J., Peters, H.A., Krasovec, M., Look, M.P., Cimerman, N., Meijer-van Gelder, M.E., Henzen-Logmans, S.C., van Putten, W.L. and Klijn, J.G. (1998) Prognostic significance of cathepsins B and L in primary human breast cancer. J. Clin. Oncol., 16, 1013–1021.

38. Clark, G.J., Kinch, M.S., Gilmer, T.M., Burridge, K. and Der, C.J. (1996) Overexpression of the Ras-related TC21/R-Ras2 protein may contribute to the development of human breast cancers. Oncogene, 12, 169–176.

39. Lee, S.W., Reimer, C.L., Fang, L., Iruela-Arispe, M.L. and Aaronson, S.A. (2000) Overexpression of kinase-associated phosphatase (KAP) in breast and prostate cancer and inhibition of the transformed phenotype by antisense KAP expression. Mol. Cell. Biol., 20, 1723–1732.

40. Maus, M.V., Reilly, S.C. and Clevenger, C.V. (1999) Prolactin as a chemoattractant for human breast carcinoma. Endocrinology, 140, 5447–5450.

41. Reynolds, C., Montone, K.T., Powell, C.M., Tomaszewski, J. E and Clevenger, C.V. (1997) Expression of prolactin and its receptor in human breast carcinoma. Endocrinology, 138, 5555–5560.

42. Bertucci, F., Van Hulst, S., Bernard, K., Loriod, B., Granjeaud, S., Tagett, R., Starkey, M., Nguyen, C., Jordan, B. and Birnbaum, D. (1999) Expression scanning of an array of growth control genes in human tumor cell lines. Oncogene, 18, 3905–3912.

43. Theillet, C., Adelaide, J., Louason, G., Bonnet-Dorion, F., Jacquemier, J., Adnane, J., Longy, M., Katsaros, D., Sismondi, P., Gaudray, P. et al. (1993) FGFRI and PLAT genes and DNA amplification at 8p12 in breast and ovarian cancers. Genes Chromosomes Cancer, 7, 219–226.

44. Bertucci, F., Bernard, K., Loriod, B., Chang, Y.C., Granjeaud, S., Birnbaum, D., Nguyen, C., Peck, K. and Jordan, B.R. (1999) Sensitivity issues in DNA array-based expression measurements and performance of nylon micro-arrays for small samples. Hum. Mol. Genet., 8, 1715–1722.

45. Eisen, M.B., Spellman, P.T., Brown, P.O. and Botstein, D. (1998) Cluster analysis and display of genome-wide expression patterns. Proc. Natl Acad. Sci. USA, 95, 14863–14868.

46. Kaplan, E. and Meier, P. (1958) Nonparametric estimation from incomplete observations. J. Am. Stat. Assoc., 53, 457–481.