



**HAL**  
open science

## **A 38-gene expression signature to predict metastasis risk in node-positive breast cancer after systemic adjuvant chemotherapy: a genomic substudy of PACS01 clinical trial**

Pascal Jézéquel, Mario Campone, Henri H. Roché, Wilfried Gouraud, Catherine Charbonnel, Gabriel Ricolleau, Florence Magrangeas, Stéphane Minvielle, Jean Genève, Anne-Laure Martin, et al.

### **► To cite this version:**

Pascal Jézéquel, Mario Campone, Henri H. Roché, Wilfried Gouraud, Catherine Charbonnel, et al.. A 38-gene expression signature to predict metastasis risk in node-positive breast cancer after systemic adjuvant chemotherapy: a genomic substudy of PACS01 clinical trial. *Breast Cancer Research and Treatment*, 2008, 116 (3), pp.509-520. 10.1007/s10549-008-0250-8 . hal-00478287

**HAL Id: hal-00478287**

**<https://hal.science/hal-00478287>**

Submitted on 30 Apr 2010

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

## A 38-gene expression signature to predict metastasis risk in node-positive breast cancer after systemic adjuvant chemotherapy: a genomic substudy of PACS01 clinical trial

Pascal Jézéquel · Mario Campone · Henri Roché · Wilfried Gouraud · Catherine Charbonnel · Gabriel Ricolleau · Florence Magrangeas · Stéphane Minvielle · Jean Genève · Anne-Laure Martin · Régis Bataille · Loïc Campion

Received: 7 November 2008 / Accepted: 7 November 2008 / Published online: 20 November 2008  
© Springer Science+Business Media, LLC. 2008

**Abstract** Currently, no prognostic gene-expression signature (GES) established from node-positive breast cancer cohorts, able to predict evolution after systemic adjuvant chemotherapy, exists. Gene-expression profiles of 252 node-positive breast cancer patients (median follow-up: 7.7 years), mostly included in a randomized clinical trial (PACS01), receiving systemic adjuvant regimen, were determined by means of cDNA custom array. In the training cohort, we established a GES composed of 38 genes (38-GES) for the purpose of predicting metastasis-free survival. The 38-GES yielded unadjusted hazard ratio (HR) of 4.86 (95% confidence interval = 2.76–8.56). Even when adjusted with the best two clinicopathological prognostic

indexes: Nottingham prognostic index (NPI) and Adjuvant!, 38-GES HRs were 3.30 (1.81–5.99) and 3.40 (1.85–6.24), respectively. Furthermore, 38-GES improved NPI and Adjuvant! classification. In particular, NPI intermediate-risk patients were divided into 2/3 close to low-risk group and 1/3 close to high-risk group (HR = 6.97 [2.51–19.36]). Similarly, Adjuvant! intermediate-risk patients were divided into 2/3 close to low-risk group and 1/3 close to high-risk group (HR = 4.34 [1.64–11.48]). The 38-GES was validated on gene-expression datasets from three external node-positive breast cancer subcohorts ( $n = 224$ ) generated from different microarray platforms, with HR = 2.95 (1.74–5.01). Moreover, 38-GES showed prognostic performance in supplementary cohorts with different lymph-node status and endpoints (1,040 new patients). The 38-GES represents a robust tool able to type systemic adjuvant treated node-positive patients at high risk of metastatic relapse, and is

Pascal Jézéquel and Loïc Campion contributed equally to this work.

**Electronic supplementary material** The online version of this article (doi:10.1007/s10549-008-0250-8) contains supplementary material, which is available to authorized users.

P. Jézéquel · W. Gouraud · C. Charbonnel · F. Magrangeas · S. Minvielle  
Unité Mixte de Génomique du Cancer, Hôpital Laënnec,  
Bd J. Monod, 44805 Nantes, Saint Herblain Cedex, France

P. Jézéquel (✉) · G. Ricolleau · R. Bataille  
Département de Biologie Oncologique, Centre de Lutte Contre le  
Cancer René Gauducheau, Bd J. Monod, 44805 Nantes,  
Saint Herblain Cedex, France  
e-mail: p-jezequel@nantes.fnclcc.fr

M. Campone  
Service d'Oncologie Médicale, Centre de Lutte Contre le Cancer  
René Gauducheau, Bd J. Monod, 44805 Nantes, Saint Herblain  
Cedex, France

H. Roché  
Service d'Oncologie Médicale, Institut Claudius Regaud, 20-24  
rue du Pont Saint Pierre, 31052 Toulouse Cedex, France

W. Gouraud · C. Charbonnel · F. Magrangeas · S. Minvielle ·  
L. Campion  
INSERM U892, Institut de Biologie, 9 quai Moncoussu,  
44035 Nantes Cedex, France

W. Gouraud · C. Charbonnel · L. Campion  
Unité de Biostatistique, Centre de Lutte Contre le Cancer René  
Gauducheau, Bd J. Monod, 44805 Nantes, Saint Herblain Cedex,  
France

J. Genève · A.-L. Martin  
Fédération Nationale des Centres de Lutte Contre le Cancer  
(FNCLCC), 101 rue de Tolbiac, 75654 Paris Cedex 13, France

especially powerful to refine NPI and Adjuvant! classification for those patients.

**Keywords** Adjuvant chemotherapy · Breast cancer · Clinicogenomic model · FEC100 regimen · Genomic study · PACS01

## Introduction

In breast cancer, gene expression profiling permitted first to distinguish breast tumors into subclasses with clinical implications, and second to predict evolution of node-negative breast cancer patients [1–5]. The initial focus on node-negative breast cancer subtyping resulted from the crucial decision of not treating low-risk patients who might have a favorable evolution without systemic adjuvant chemotherapy. Currently, three prognostic breast cancer tests resulting from gene expression studies are commercially available: Oncotype Dx<sup>®</sup> (Genomic Health, Redwood City, CA, USA), MammaPrint<sup>®</sup> (Agendia BV, Amsterdam, The Netherlands) and H/I<sup>®</sup> (AvariaDX, Carlsbad, CA, USA) [3, 6, 7]. All these tests are designed to be used in the management of node-negative breast cancer. Today, in breast cancer, no robust gene-expression signature (GES) based exclusively on node-positive patients treated by systemic adjuvant chemotherapy, and able to predict metastasis evolution, exists. Being able to type those node-positive patients would lead to orientate and include high-risk patients in new clinical trials testing novel targeted therapies. The goal of our study was thus to identify and validate a GES for the purpose of predicting metastatic relapse in node-positive patients treated by systemic adjuvant chemotherapy, and so, to enlarge the panel of prediction of breast cancer GESs.

## Methods

### Patients

The training cohort (TC) was composed of 252 women who had unilateral breast cancer, were node-positive, showed no evidence of metastasis at diagnosis, were younger than 65 years of age, and treated by adjuvant systemic chemotherapy. The majority of the studied patients (90.5%) was primarily included in a multicentric phase III clinical trial (PACS01) conducted by investigators from the French Federation of Cancer Center [8]. One hundred and twenty eight patients were also part of a previous study; in the present work, follow-up data were actualized [9]. Twenty four patients (9.5%) followed in the René Gauducheau Cancer Center were included based on the same inclusion

criteria as those used for PACS01 trial. Briefly, after surgery, patients received intravenous adjuvant treatment with 5-fluorouracil 500 mg/m<sup>2</sup>, epirubicin 100 mg/m<sup>2</sup> and cyclophosphamide 500 mg/m<sup>2</sup> (FEC100) every 21 days for 6 cycles (PACS01-FEC arm and René Gauducheau Cancer Center), or 3 cycles (PACS01-FEC+Docetaxel arm). Patients were followed up for metastasis-free survival (MFS). No bioclinical or MFS heterogeneity between the 3 subpopulations was found, allowing their pooling (supplementary Table 1). The median follow-up was 7.7 years (range 1.22–10.02). During the period, 65 patients showed evidence of distant MR and 47 died. This substudy of PACS01 trial was reviewed and approved by the ethics committee/institutional (CCPPRB) review board (number 1-97-13). All patients signed informed consent for research purpose.

External validation was based on independent breast cancer patients with the following inclusion criteria: available genomic profiles, node-positive and metastatic relapse (MR) well defined in order to use the same outcome definition as the one used in the TC. Hence, we extracted three subcohorts (validation cohorts [VCs] 1–3) from three published cohorts (Table 1, supplementary Table 2) [4, 10, 11]. No statistical heterogeneity was found between TC and VC1 + VC2 + VC3 for available data (estrogen receptor [ER] and 5-year MFS) (supplementary Table 3).

**Table 1** Methodological distribution and status of the studied cohorts

Code	Data sets	Nodal status			Clinical outcome	
		N+	N–	Mixed	MR	ue
Training cohort						
TC	Our study	252			65	
	Total	252			65	
Validation cohorts						
VC1	Van de Vijver et al. [4]	144			47	
VC2	Sotiriou et al. [10]	53			23	
VC3	Sotiriou et al. [11]	27			7	
	Total	224			77	
Exploratory cohorts						
EC1	Van de Vijver et al. [4]	151			54	
EC2	Sotiriou et al. [10]	46			7	
EC3	Sotiriou et al. [11]	149			33	
EC4	Wang et al. [5]	286			107	
	Subtotal	632			201	
EC5	Pawitan et al. [12]		159			40
EC6	Ivshina et al. [13]		249			89
	Subtotal		408			129

MR metastatic relapse; ue unspecified event: local relapse, metastatic relapse or cancer death; TC training cohort; VC validation cohort 1–3; EC, exploratory cohort 1–6

To explore performance capacities of our gene signature (*n*-GES), a complementary study was conducted on breast cancer patients (exploratory cohorts [ECs] 1–6) who did not fulfill *n*-GES establishment inclusion criteria (node-mixed, node-negative and unspecified event [local relapse, metastatic relapse or cancer death]) (Table 1, supplementary Table 2) [4, 5, 10–13].

#### Tumor tissue samples and RNA isolation

All tumor tissue samples were surgically collected and immediately macrodissected by a pathologist, snap-frozen in liquid nitrogen and stored until RNA extraction. Total RNA was prepared by the CsCl-cushion, as described [14]. The quality and the quantity of RNA samples were evaluated by Agilent 2100 Bioanalyser RNA LabChip kit (Agilent Technologies, Palo Alto, CA, USA).

#### Nylon microarray technology

We used nylon microarrays manufactured in our lab, which contained 8,032 unique sequenced-verified cDNA clones, representing 5,776 distinct genes, chosen using the expressed sequence tag database from the NCBI with update Hs#196 (<http://www.ncbi.nlm.nih.gov>). The same microarray tools were used and described in previous works [9, 15, 16]. Microarray characteristics and data have been deposited in the NIH Gene Expression Omnibus (Series accession number: GSE11264) according to minimum information about a microarray experiment (<http://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE11264>).

#### Statistical analysis

##### *Bioinformatics*

DNA microarrays were scanned at 25- $\mu$ m resolution using a Fuji BAS 5000 image plate system (Raytest, Paris, France). The hybridization signals were quantified using ArrayGauge software v.1.3 (Fuji Ltd, Tokyo, Japan). For each membrane, the raw data were background noise subtracted (median value of negative controls  $\pm$  6 standard deviation), then normalized by global intensity of hybridization. Patient data were adjusted by the amount of polymerase chain reaction product spotted onto the membrane,  $\log_2$ -transformed and finally standardized (mean = 0; standard deviation = 1). Only genes well measured for at least 70% of patients were retained.

##### *Establishment of the gene expression signature*

Genes were selected by means of resampling and univariate Cox permutations techniques to minimize overfitting and to

maximize stability of gene list for MFS prediction. We created 100 different random samples composed of the 2/3 of TC. For each of these 100 subsamples  $i$  ( $i = 1-100$ ), we determined an unbiased predictor ( $p_i$ ) with internal robust cross-validation (leave-one-out cross-validation method). One-hundred  $p_i$  ( $i = 1-100$ ) were thus determined. In order to increase stability of final gene list ( $n$ ), only genes present in at least 50% of the unbiased resampled gene predictors  $p_i$  were retained to build the *n*-GES. The *n*-GES score was calculated as the difference between the mean of the genes overexpressed genes in MR patients and the mean of the genes underexpressed genes in MR patients. To limit overfitting no weight was assigned to genes. Standardization of gene distributions permitted extrapolation of fully determined *n*-GES (including cutoff) onto other cohorts. Optimal cutoff between *n*-GES low-risk and high-risk groups was determined as the value with minimal Akaike Information Criterion (AIC). The *n*-GES was evaluated in VCs, separately and pooled. When at least one GES gene was missing in an external dataset, we calculated a reduced 38-GES and determined its optimal cutoff in TC. Moreover it was necessary to evaluate if this subscore was significantly correlated with 38-GES complete score to permit its use (Supplementary Appendix 1).

##### *Survival analysis and endpoints*

Time from surgery to MR (primary endpoint) was retained for GES establishment. Overall survival (OS), defined as time from surgery to death from any cause, was used to reinforce MFS analysis in TC. Survival curves were plotted according to Kaplan–Meier method and compared by means of log-rank test (or Wilcoxon log-rank test if number of events was small, in order to give greater weight to earlier events). For each univariate survival analysis, 10,000 permutation survival analyses were performed to ensure  $P$  values were not optimized by overfitting. Only parameters with permutation  $P < 0.05$  were entered into the forward Cox proportional hazards model. Proportional hazards assumption was verified for the final models by means of Schoenfeld residuals study and  $-\ln(-\ln(\text{survival}))$  curves. Because time to distant metastasis was not specified for EC5 and EC6, time to the first cancer event (local relapse, MR or death) was used.

##### *Sensitivity, specificity*

Receiver operating characteristic (ROC) analysis with MR within 5 years as a defining point was computed. The area under the curve (AUC) was used as a measure of the signature global performance in the test set. Time-dependent ROC curves were calculated for *n*-GES and the best two clinical prognostic indexes: Nottingham prognostic index (NPI) and Adjuvant! (10-year OS score was used) [17].

### Clinicogenomic model establishment and classification

Firstly, multivariate Cox regression analysis was used to determine if *n*-GES added independent prognostic information to NPI and Adjuvant!. To rule out the possibility that our results were dependent on the choice of clinical indexes cutoff, we carried out a sensitivity analysis in which we varied the NPI raw score that defined low risk, as predicted by NPI, from 3.5 (in which case most patients were classified in the high-risk group) to 7 (in which case most patients were classified in the low-risk group). For Adjuvant!, the same analysis was done from 20% (in which case most patients were classified in the low-risk group) to 80% (in which case most patients were classified in the high-risk group). Secondly, we focused on patients with NPI and Adjuvant! intermediate risk: score 2 and 10-year OS probability between 60 and 80%, respectively. For those patients, the therapeutic decision could be facilitated by an appropriate refining, for example to orientate *n*-GES high-risk patients towards new clinical trials.

All analyses were performed with SAS 9.1 (SAS Institute Ins., Cary, NC), Stata 10.0 SE (Stata Corp., College Station, TX), R software (version 2.5.1) and BRB ArrayTools developed by Dr. R. Simon and A. Peng (Bethesda, MD; 2003) (<http://linus.nci.nih.gov/BRB-ArrayTools.html>). All statistical tests were two-sided and statistical significance was defined as *P* less than 0.05.

## Results

The aim of this study was to examine whether the *n*-GES; (1) had prognostic value in node-positive patients treated by chemotherapy independently of known prognostic parameters such as age, Scarff-Bloom-Richardson grade, tumoral size, ER or number of positive nodes (alone or combined into NPI and Adjuvant!); (2) could improve classification of NPI and Adjuvant!, especially for intermediate-risk patients; (3) confirmed its performance in node-positive independent cohort; (4) had also prognostic value in node-negative patients.

### Establishment of the gene expression signature

Thirty eight genes were kept (Table 2). Control of false discovery rate (FDR) among the differentially expressed genes was assessed by 10,000 permutations: permutation *P* was <0.008 and FDR was <18% for all genes but one (*P* = 0.014, FDR = 24%) [18]. Calculation of 38-GES relapse score is defined below (Panel). NPI and Adjuvant! scores are calculated as shown in Supplementary Appendix 2. Optimal cutoff determination for 38-GES, NPI and Adjuvant! is described in supplementary Fig. 1.

### Panel: Calculation of 38-GES score

$$38\text{-GES} = \frac{1}{30} \sum_{i=1}^{30} G_i - \frac{1}{8} \sum_{j=1}^8 G_j$$

$G_i$  and  $G_j$  represent overexpressed and underexpressed gene-expression standardized values, respectively. No weight was assigned to genes in order to limit overfitting. 38-GES optimal cutoff in TC was 0.005.

### Comparison of 38-GES with NPI and Adjuvant!

The unadjusted MFS and OS hazard ratios of classical bioclinical factors and best clinicopathologic risk classifications are shown in supplementary Tables 4 and 5, respectively. The prognostic value of the gene signature was stronger than each of traditional risk factors, NPI and Adjuvant! classifications. Optimal MFS unadjusted hazard ratio (HR) (high risk vs. low risk) for 38-GES, NPI and Adjuvant! were 4.86 (95% confidence interval [CI] 2.76–8.56; *P* < 0.0001), 4.48 (95% CI 2.62–7.66; *P* < 0.0001), 4.16 (95% CI 2.54–6.80; *P* < 0.0001), respectively (supplementary Table 4). AIC of the 3 classification schemes showed that the classification power of 38-GES (AIC = 641.44) was better than the NPI's one (AIC = 649.58), which itself was better than Adjuvant!'s one (AIC = 653.79) to explain the MFS. ROC analysis for MR within 5 years, as a defining point, showed that 38-GES global predictive performance (AUC = 0.779) was higher than NPI's one (AUC = 0.757) and Adjuvant!'s one (AUC = 0.751) with higher sensitivity (Fig. 1). Time-dependent ROC curves confirmed that the global predictive performance of the 38-GES was better than that of NPI and Adjuvant! (data not shown).

### Independent prognostic value of 38-GES

The gene signature adjusted hazard ratios for NPI and for the 10-year OS probability Adjuvant! score were 3.30 (95% CI 1.81–5.99) and 3.40 (95% CI 1.85–6.24), respectively (supplementary Table 4). Moreover sensitivity analysis showed that the 38-GES kept its independent prognostic value whatever the definition of NPI and Adjuvant! high-risk group was (supplementary Fig. 2). These results were also confirmed for OS analysis (supplementary Table 5, supplementary Fig. 3). So, relevance of a clinicogenomic model combining 38-GES and bioclinical parameters was shown.

### Clinicogenomic model establishment and classification

38-GES refined NPI and Adjuvant! classification, in particular for intermediate-risk patients (Fig. 2, supplementary

**Table 2** The 38-gene expression signature for predicting breast cancer metastatic relapse

% Model <i>n</i> = 100	Permutation <i>P</i> <i>n</i> = 10,000	Gene symbol	Gene title	MR sense	Biological process (Gene ontology)	Genbank #
99	<1.10 <sup>-7</sup>	<i>BTG2</i>	BTG family, member 2	–	DNA repair, transcription, regulation of transcription DNA-dependent, negative regulation of cell proliferation	NM_006763
96	<1.10 <sup>-7</sup>	<i>UBE2C</i>	Ubiquitin-conjugating enzyme E2C	+	Ubiquitin proteasome system, cell cycle	BC032677
96	0.0001	<i>VEGFA</i>	Vascular endothelial growth factor A	+	Angiogenesis, regulation of progression, signal transduction	NM_001025366
94	0.0041	<i>SLC25A5</i>	Solute carrier family 25, member 5	+	Transport, mitochondrial transport	R56229
93	0.0010	<i>FABP5</i>	Fatty acid binding protein 5	+	Lipid metabolic process, transport	BG282526
89	0.0004	<i>ENO1</i>	Enolase 1, (alpha)	+	Glycolysis, transcription	AL833741
86	0.0004	<i>PSMB7</i>	Proteasome (prosome, macropain) subunit, beta type, 7	+	Ubiquitin proteasome system, cell cycle	R54562
81	0.0003	<i>TUBB6</i>	Tubulin, beta 6	+	Microtubule-based process, microtubule-based movement, protein polymerization	NM_032525
81	0.0016	<i>TACC3</i>	Transforming, acidic coiled-coil containing protein 3	+	Progression, response to stress, hemopoiesis	BC106071
80	0.0001	<i>SURF4</i>	Surfeit 4	+	<i>Unknown</i>	NM_033161
77	0.0007	<i>RHOC</i>	Ras homolog gene family, member C	+	Glycolysis, signal transduction, regulation NF-kappaB cascade	AK094474
76	0.0004	<i>CHCHD2</i>	Coiled-coil-helix- coiled-coil-helix domain containing 2	+	<i>Unknown</i>	NM_016139
71	0.0011	<i>SFRS5</i>	Splicing factor, arginine/serine- rich 5	–	mRNA splice site selection, mRNA processing, RNA splicing	R28509
70	0.0004	<i>TNFSF13</i>	Tumor necrosis factor (ligand) superfamily, member 13	–	Angiogenesis, cell proliferation, apoptosis, immune response	NM_003808
70	0.0024	<i>DDT</i>	D-dopachrome tautomerase	+	Melanin biosynthetic process from tyrosin	AI936198
67	0.0003	<i>C1orf64</i>	Chromosome 1 open reading frame 64	–	<i>Unknown</i>	AI732325
67	0.0040	<i>RAB10</i>	RAB10, member RAS oncogene family	+	Regulation of transcription DNA-dependent, intracellular protein transport, nucleocytoplasmic transport, signal transduction, small GTPase mediated signal transduction	AA074077
66	0.0011	<i>XIST</i>	X (inactive)-specific transcript	–	<i>Unknown</i>	H78857
65	0.0017	<i>TUBA4A</i>	Tubulin, alpha 4a	+	Microtubule-based process, microtubule-based movement, protein polymerization	AK054731
63	0.0041	<i>PPA1</i>	Pyrophosphatase (inorganic) 1	+	Phosphate metabolic process	BF694769
62	0.0011	<i>C4B</i>	Complement component 4 B	–	Inflammatory response, complement activation, innate immune response	NM_001002029
62	0.0018	<i>TUBB2C</i>	Tubulin, beta 2C	+	Cell motility, microtubule-based process, microtubule- based movement, natural killer cell mediated cytotoxicity, protein polymerization	BX648521
62	0.0071	<i>C3orf37</i>	Chromosome 3 open reading frame 37	+	<i>Unknown</i>	T65434

**Table 2** continued

% Model	Permutation $P$ $n = 10,000$ $n = 100$	Gene symbol	Gene title	MR sense	Biological process (Gene ontology)	Genbank #
60	0.0017	<i>GATA3</i>	GATA binding protein 3	–	Cell fate determination, transcription, regulation of transcription DNA-dependent, transcription from RNA polymerase II promoter, defense response	R16918
60	0.0039	<i>PGK1</i>	Phosphoglycerate kinase 1	+	Glycolysis, phosphorylation	H17787
59	0.0014	<i>PPP1CA</i>	Protein phosphatase 1, catalytic subunit, alpha isoform	+	Carbohydrate metabolic process, protein amino-acid dephosphorylation, cell cycle	NM_001008709
59	0.0022	<i>BCL2A1</i>	BCL2-related protein A1	+	Anti-apoptosis, regulation of apoptosis	N28416
58	0.0005	<i>IL6ST</i>	Interleukin 6 signal transducer (gp130, oncostatin M receptor)	–	Glycogen metabolic process, immune response, signal transduction, cell surface receptor linked signal transduction, positive regulation of cell proliferation	H04779
58	0.0017	<i>PSMB3</i>	Proteasome (prosome, macropain) subunit, beta type, 3	+	Ubiquitin-dependent protein catabolic process	R00911
58	0.0020	<i>RPN1</i>	Ribophorin I	+	Protein amino-acid glycosylation	NM_002950
57	0.0033	<i>ETFB</i>	Electron-transfer-flavoprotein, beta polypeptide	+	Electron transport, transport	R55473
56	0.0006	<i>CYC1</i>	Cytochrome c-1	+	Electron transport, transport	BF569085
56	0.0144	<i>POR</i>	P450 (cytochrome) oxidoreductase	+	Electron transport	CD014011
55	0.0016	<i>TUBB3</i>	Tubulin, beta 3	+	Microtubule-based process, microtubule-based movement, mitosis, signal transduction, G-protein coupled receptor protein signaling pathway	CR596505
53	0.0060	<i>TMED9</i>	Transmembrane emp24 protein transport containing 9	+	Transport	NM_017510
52	0.0036	<i>TOE1</i>	Target of EGR1, member 1 (nuclear)	+	Negative regulation of transcription from RNA polymerase II promoter	H29584
52	0.0045	<i>C17orf37</i>	Chromosome 17 open reading frame 37	+	Cell redox homeostasis, protein homooligomerization	T84927
51	0.0073	<i>PSMA5</i>	Proteasome (prosome, macropain) subunit, alpha type, 5	+	Ubiquitin-dependent protein catabolic process	H48425

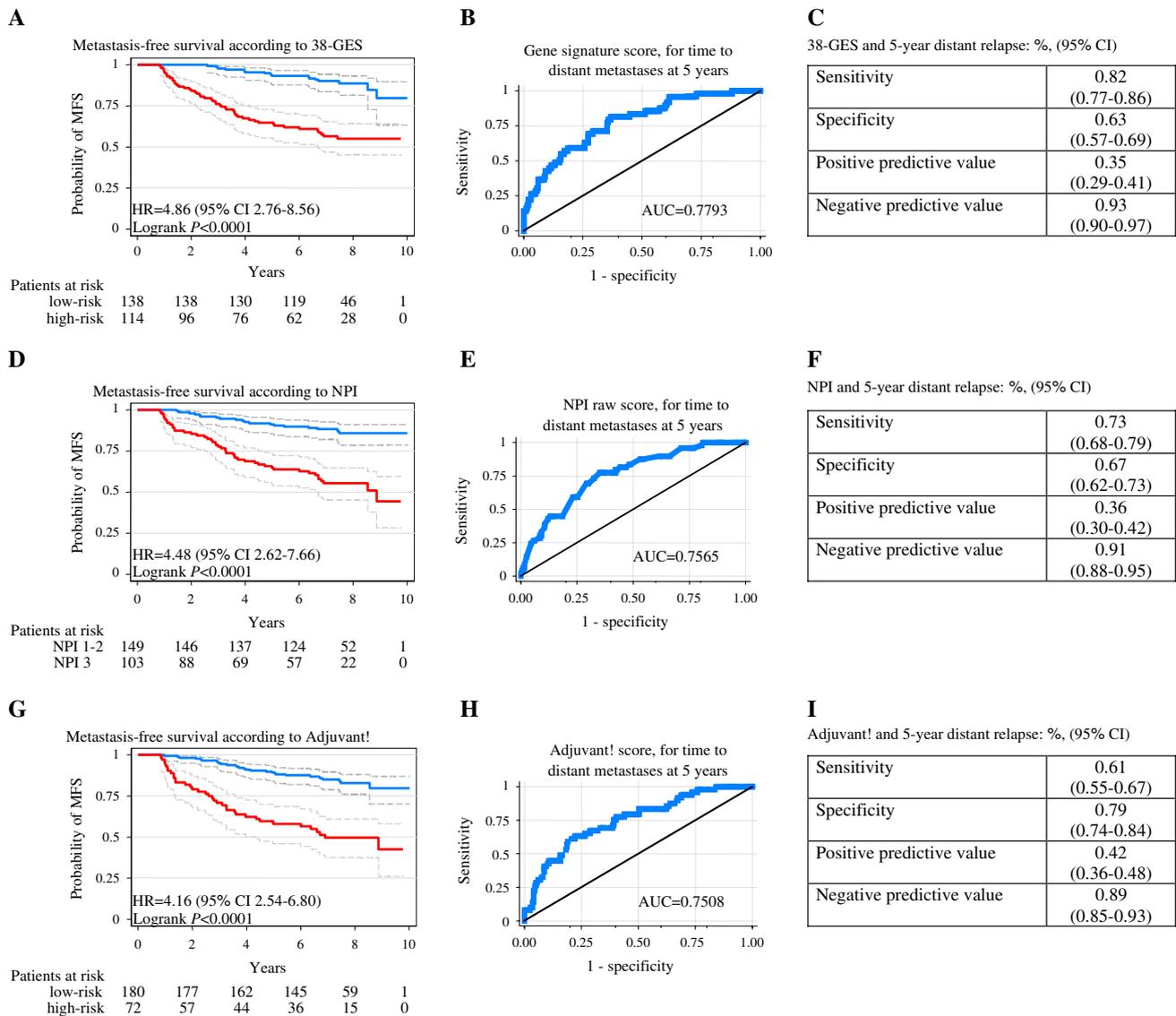
MR metastatic relapse

Table 6). NPI 2 patients with 5-year MFS = 90.4% could be separated in 2/3 with 5-year MFS = 96.7% close to NPI 1 and 1/3 5-year MFS = 76.7% close to NPI 3 (HR = 6.97,  $P < 0.0001$ ). Similarly, Adjuvant! 60–79% 10-year OS patients with 5-year MFS = 88.9% could be separated in 2/3 with 5-year MFS = 93.7% close to Adjuvant!  $\geq 80\%$  10-year OS group and 1/3 5-year MFS = 76.9% close to Adjuvant!  $< 60\%$  10-year OS group (HR = 4.34,  $P < 0.0001$ ). The same successful stratification was obtained for NPI and Adjuvant! high-risk patients (supplementary Table 6, supplementary Fig. 4).

Validation of the 38-GES on independent node-positive patients

Correlation study permitted us to apply reduced 38-GES on external datasets with missing genes. Coefficients range was 0.973–0.989 with all  $P < 0.00001$  (Supplementary Appendix 1).

38-GES significantly predicted metastatic evolution in pooled node-positive VC (HR = 2.95, 95% CI 1.74–5.01,  $P < 0.0001$ ) with good sensitivity (0.84) and good global predictive performance (AUC = 0.73) (Fig. 3). ER was the



**Fig. 1** Kaplan–Meier analysis for metastasis-free survival, receiver operating characteristic curves and characteristics details in training cohort. **a–c** 38-GES. **d–f** Nottingham prognostic index. **g–i** Adjuvant! scores

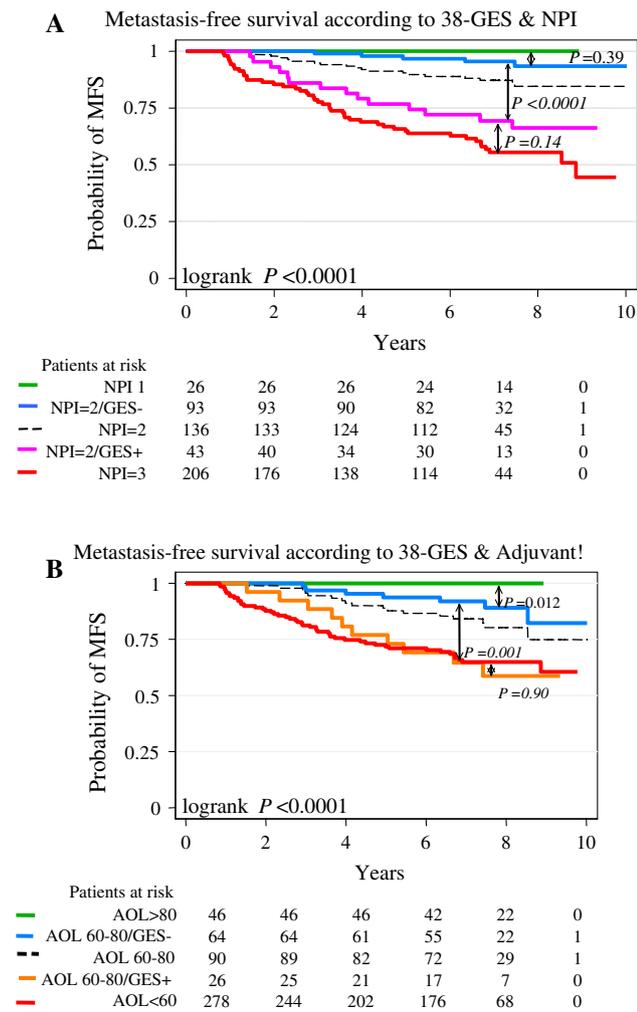
only parameter which was available in each VC. Even when ER adjusted, 38-GES kept strong prognostic value in pooled VC (HR = 2.63, 95% CI 1.50–4.62,  $P < 0.0001$ ) and separate VCs (Fig. 4).

Assessment of the 38-GES performance on node-negative and “node-mixed” cohorts

38-GES significantly predicted MFS in pooled node-negative ECs (ER adjusted HR = 2.11, 95% CI 1.53–2.92,  $P < 0.0001$ ) and time to disease event (local relapse, metastasis and death) in mixed node cohorts (ER adjusted HR = 2.77, 95% CI 1.93–3.97,  $P < 0.0001$ ) (Fig. 4 and supplementary Fig. 5).

Molecular biology of the markers

A large majority (86%) of the 38 genes showed a significant Cox univariate  $P$ -value in at least one external subcohort or cohort (Table 3). Genes reaching or exceeding 50% of significance in the 10 groups are: *UBE2C*, *TACC3*, *IL6ST*, *TUBB2C*, *TUBB3*, *PGK1*, *PPP1CA* and *SLC25A5*. Results displayed in Table 3 and literature mining confirmed that the vast majority of the genes retained in our 38-GES are biologically relevant and linked to breast cancer. Several biological processes were identified by using Gene Ontology database (<http://www.geneontology.org/>). In order of representation, these were: transcription (6 genes), transport (6 genes), ubiquitin-proteasome system (4 genes),



**Fig. 2** Kaplan–Meier analysis for clinicogenomic models. **a** 38-GES and NPI. **b** 38-GES and Adjuvant!. 38-GES stratifies intermediate-risk patients based on NPI and Adjuvant! classifications

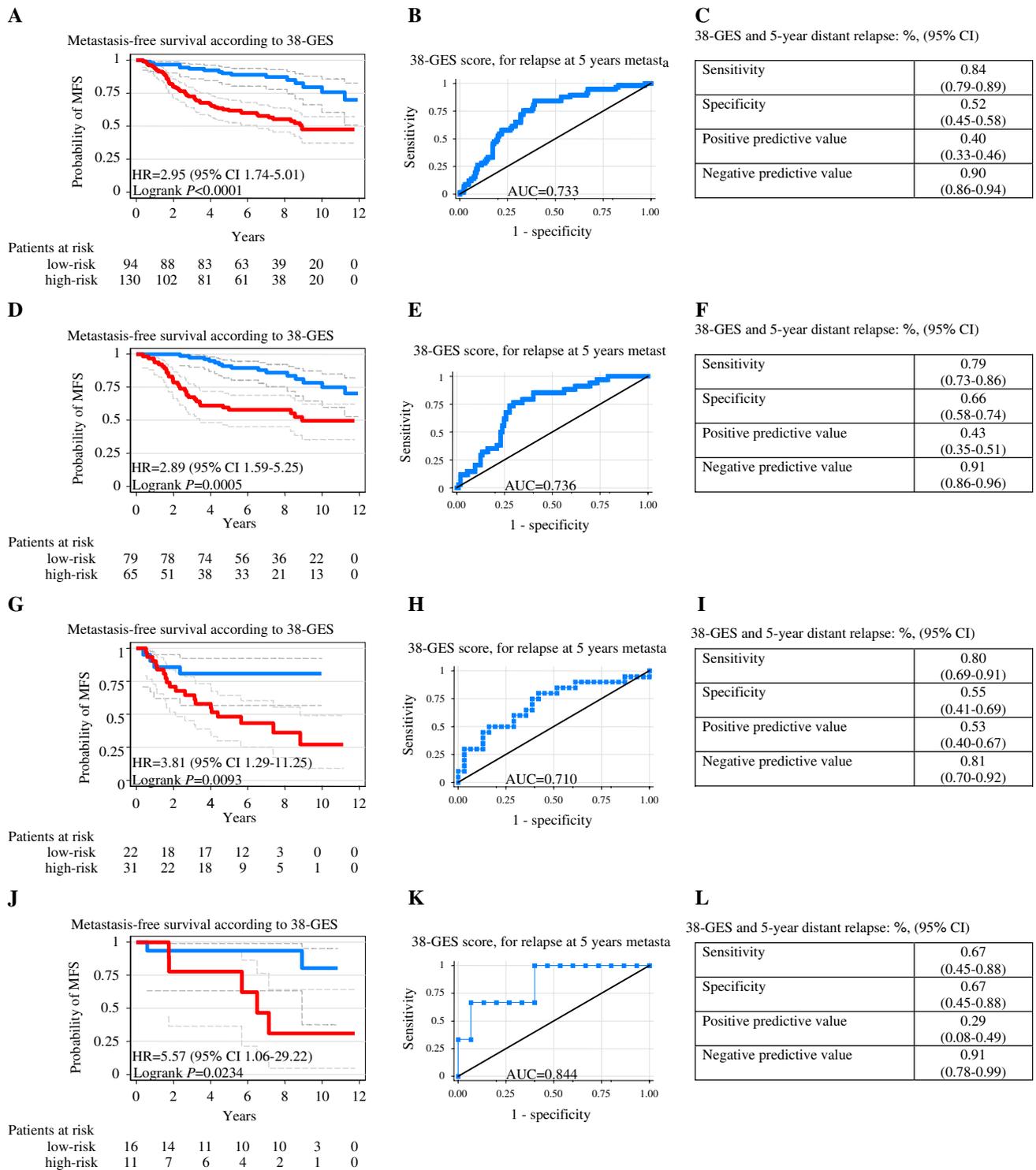
microtubule-based process (4 genes), signal transduction (4 genes), proliferation (3 genes), cell cycle (3 genes), glycolysis (3 genes), immune response (3 genes), electron transport (3 genes), apoptosis (2 genes), angiogenesis (2 genes) and progression (2 genes).

## Discussion

The strengths of our study are detailed below. Our methodology was based on rigorous internal cross-validation for establishment of a GES on a training cohort (representative of node-positive population), followed by external validation using clinically relevant patients. Furthermore, the size of our training cohort ( $n = 252$ ) limited overoptimistic performance estimation due to model overfit and multiple training/test partitions maximized the stability of the obtained 38-gene list. This cohort was homogeneous,

permitting identification of a GES not subject to potential confounding factors related to lymph node status or different systemic adjuvant treatment. Training and validation cohort homogeneity was verified and permitted extrapolation. Prognostic value of 38-GES was tested in pooled validation cohorts but also in each of validation subcohorts separately. As recommended for validation, predictor was fully developed from training cohort (including cutoff value), and predictive performance was verified using the same outcome definition. Our 38-GES was more powerful than the best classical prognostic references (NPI and Adjuvant!) and added independent prognostic information [19, 20]. Importantly, the prognostic performance of our 38-GES was confirmed in several independent data sets generated from different microarray platforms (Affymetrix, Agilent, academic cDNA [NCI]). Agreement study of the 38-GES gene expression significativities showed concordant results, which represents a response to skeptical view of microarray studies (Table 3). These results confirmed on one hand the preanalytical and analytical robustness (different cohorts, biobanks and genomic protocols [RNA isolation, probe preparation and labeling, hybridization, microarray platforms]) of the vast majority of the 38 selected genes, and on the other hand their critical role in breast cancer prediction.

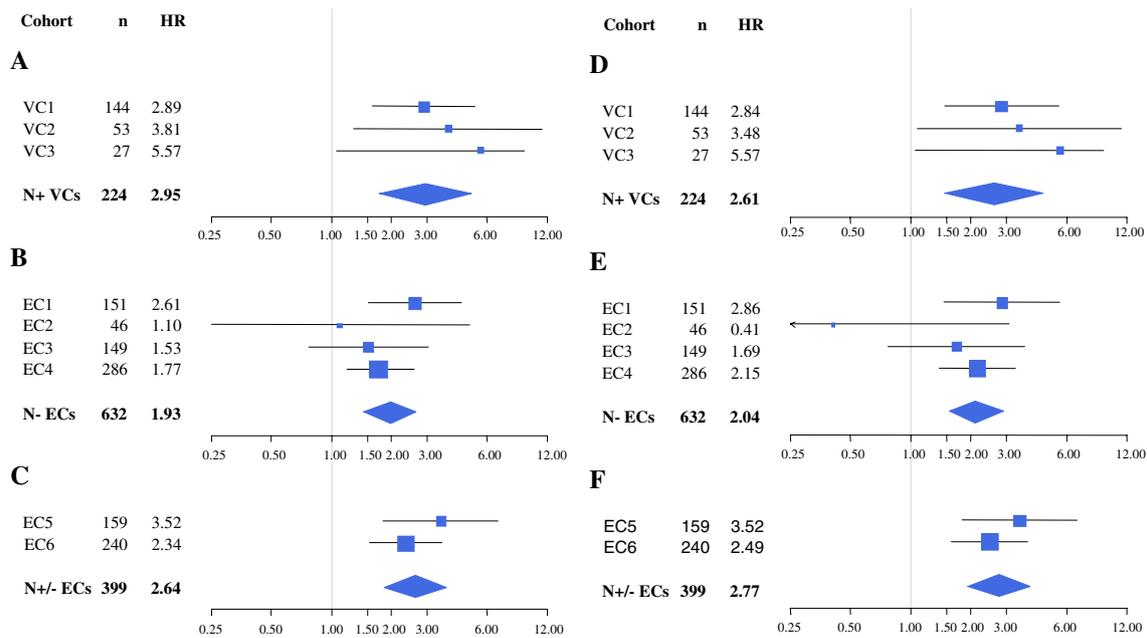
Three main results therefore emerged from this study. First, the 38-GES represents a robust tool able to type systemic adjuvant treated node-positive patients at high risk of MR, and is especially powerful to refine NPI and Adjuvant! classifications. Hence, 38-GES high-risk patients should be orientated and included in new clinical trials testing new targeted therapies, while low-risk patients could still receive node-positive standard treatment. Second, genes, or cellular pathways, revealed by 38-GES might be used as targets for specific treatments of breast cancer. Third, after proper comparison with Mammaprint<sup>®</sup> and Oncotype Dx<sup>®</sup>, 38-GES could be used as a prognostic tools for node-negative breast cancer patients. Actually, since their initial description, prediction abilities of the two main GES in node-negative breast cancer (Mammaprint<sup>®</sup> [70-GES] and Oncotype Dx<sup>®</sup> [21-GES]) have evolved due to complementary studies. The 70-GES was developed for prognostic prediction of untreated node-negative breast cancer patients [3]. Soon after, this GES showed a prognostic value in node-positive patients who received adjuvant systemic treatment [4]. On the contrary, the Oncotype Dx assay was developed on a node-mixed cohort to predict distant relapse and validated in positive ER node-negative breast cancer patients treated with adjuvant tamoxifen [6]. Since then, prognostic and predictive performance of this test has been refined by using different breast cancer cohorts [21–27]. In this work, we observed a comparable enlargement of the prediction spectrum of our



**Fig. 3** Kaplan-Meier analysis for metastasis-free survival, receiver operating characteristic curves and characteristics details for pooled and separate node-positive validation cohorts (VCs). **a-c** Pooled VC1, VC2 and VC3. **d-f** VC1. **g-i** VC2. **j-l** VC3

GES. In an exploratory study, we evaluated its performance value on different and various breast cancer subcohorts and cohorts. These results show that the prediction ability of the 38-GES is not restricted to MR in

FEC100 treated node-positive breast cancer patients. What can explain this disparate and unexpected spectrum of prediction? The fact that our 38-GES bears predictive and prognostic performance is not so surprising. A predictor of



**Fig. 4** Forest plots of 38-GES metastasis-free hazard ratios (HRs) and 95% confidence intervals (CIs) in high-risk vs low-risk groups for validation cohorts (VCs) and exploratory cohorts (ECs). **a–c**

Unadjusted HR. **d–f** Estrogen receptor (ER) adjusted hazard ratio. (**a, d**) VC1 to VC3. (**b, e**) EC1 to EC4. (**c, f**) EC5 and EC6

**Table 3** Classification of the 38-GES genes based on number of significant Cox univariate analysis *P* values in the different studied breast cancer cohorts

No of cohorts with significant expression	Gene	TC	VC1	VC2	VC3	EC1	EC2	EC3	EC4	EC5	EC6
7	<i>UBE2C</i>	0.00005	0.00113	NS	NS <sup>a</sup>	0.00024	NS	0.01962	0.00066	0.00051	0.00025
6	<i>TACC3</i>	0.00127	0.02032	NS	NS	0.00033	NS	NS <sup>b</sup>	0.03920	0.00100	0.00018
5	<i>IL6ST</i>	0.00034	0.01533	NS	NS <sup>b</sup>	0.03427	NS	NS	0.00812	0.00076	NS
5	<i>TUBB2C</i>	0.00096	NS <sup>a</sup>	md	NS	0.00929	md	0.04814	NS	0.00003	0.00101
5	<i>TUBB3</i>	0.00146	NS	md	NS	NS	md	0.04956	0.03492	<0.00001	0.00083
5	<i>PGK1</i>	0.00322	0.00048	md	NS	0.00006	md	NS	NS	0.00005	0.02565
5	<i>PPP1CA</i>	0.00109	0.00349	NS	NS	0.01239	NS	NS	NS	0.00471	0.00024
5	<i>SLC25A5</i>	0.00191	NS <sup>a</sup>	0.01349	NS	0.02281	NS <sup>a</sup>	NS	NS	0.00039	0.00382
4	<i>BTG2</i>	0.00001	0.00543	NS <sup>a</sup>	NS <sup>a,b</sup>	0.00020	NS <sup>a</sup>	NS <sup>a</sup>	0.01687	NS	NS
4	<i>C4B</i>	0.00061	NS <sup>a</sup>	NS	0.01467	NS <sup>a</sup>	NS	NS	NS	0.00401	0.02725
4	<i>VEGFA</i>	0.00020	NS <sup>a,b</sup>	0.01384	NS	0.00294	NS <sup>a</sup>	NS	NS	NS <sup>b</sup>	0.04090
4	<i>SFRS5</i>	0.00077	NS <sup>a</sup>	NS	NS <sup>b</sup>	0.01678	NS	NS	0.03148	0.02179	NS
4	<i>ENO1</i>	0.00008	0.01194	md	NS	0.01341	md	NS	NS	0.02788	NS
4	<i>GATA3</i>	0.00060	NS <sup>a</sup>	0.02185	NS	0.02143	NS <sup>a</sup>	NS	NS	0.01549	NS
4	<i>TMED9</i>	0.00415	0.00025	NS	NS	NS <sup>a</sup>	NS	NS	0.00002	NS	0.00192
4	<i>POR</i>	0.01509	NS <sup>a</sup>	NS	NS	0.00959	NS	NS	NS	0.00001	0.03871
3	<i>PSMB3</i>	0.00177	NS <sup>a,b</sup>	NS	NS	NS <sup>a</sup>	NS	NS	NS	0.03941	0.00096
3	<i>CYC1</i>	0.00091	0.00148	NS	NS	NS <sup>a,b</sup>	NS <sup>b</sup>	NS <sup>b</sup>	NS	0.00105	NS <sup>b</sup>
3	<i>XIST</i>	0.00027	md	md	NS <sup>b</sup>	md	md	NS	NS	0.00913	0.00941
3	<i>C17orf37</i>	0.00272	md	md	md	md	md	md	md	0.03881	0.00073
2	<i>FABP5</i>	0.00058	NS	NS <sup>a,b</sup>	NS	NS	NS <sup>a</sup>	NS	0.01527	NS <sup>b</sup>	NS
2	<i>PSMB7</i>	0.00059	NS <sup>a,b</sup>	NS	NS	NS <sup>a</sup>	NS	0.04657	NS	NS	NS
2	<i>TUBA4A</i>	0.00153	NS	NS <sup>a,b</sup>	NS	NS	NS <sup>a</sup>	NS	NS	0.01089	NS

**Table 3** continued

No of cohorts with significant expression	Gene	TC	VC1	VC2	VC3	EC1	EC2	EC3	EC4	EC5	EC6
2	<i>DDT</i>	0.00124	0.04183	NS	NS	NS <sup>a</sup>	NS	NS	NS	NS <sup>b</sup>	NS <sup>b</sup>
2	<i>CHCHD2</i>	0.00035	md	md	NS	md	md	NS	NS	0.00016	NS <sup>b</sup>
2	<i>RAB10</i>	0.00218	NS <sup>b</sup>	NS	md	NS	NS	md	md	0.00065	NS
2	<i>RPNI</i>	0.00161	NS	md	NS	NS	md	NS	NS	NS <sup>b</sup>	0.00856
2	<i>ETFB</i>	0.00175	NS	0.01651	NS	NS	NS <sup>a</sup>	NS	NS	NS	NS <sup>b</sup>
2	<i>SURF4</i>	0.00018	md	md	md	md	md	md	md	0.00051	NS
2	<i>TNFSF13</i>	0.00013	NS	NS	NS	NS	NS	NS	0.00979	NS	NS
2	<i>PPA1</i>	0.00204	NS	NS	NS	NS	NS	NS	NS	NS	0.00778
2	<i>BCL2A1</i>	0.00215	NS	NS	NS	NS	NS	NS	NS	0.01237	NS
2	<i>PSMA5</i>	0.00673	NS	NS	NS	NS	NS	NS	NS	NS	0.00589
1	<i>C1orf64</i>	0.00087	md	md	md	md	md	md	md	NS	NS <sup>b</sup>
1	<i>TUBB6</i>	0.00028	md	md	NS	md	md	NS	NS	NS <sup>b</sup>	NS
1	<i>C3orf37</i>	0.00291	NS	NS	NS	NS	NS	NS	NS	NS <sup>b</sup>	NS
1	<i>TOE1</i>	0.00296	md	md	NS	md	md	NS	NS	NS	NS
1	<i>RHOC</i>	0.00047	NS	NS	NS	NS	NS	NS	NS	NS	NS
No. of genes with:	$P < 0.05$	38	10	4	1	13	0	4	9	22	17
	$P < 0.10$	38	14	6	5	14	1	6	9	28	22

TC training cohort; VC validation cohort 1 to 3; EC exploratory cohort 1 to 6; NS non significant ( $P > 0.05$ ); md missing data

<sup>a</sup>  $P < 0.05$  in  $VC_i \cup EC_i$ ;  $i = 1, 2$  or  $3$

<sup>b</sup>  $0.05 < P < 0.10$

evolution after a systemic adjuvant treatment logically might contain predictive and prognostic informativity; but to what extent? As other researchers, we think that such predictors might represent, to a large extent, a predictor of phenotype linked to prognosis [24, 28]. A last remark interests, once again, the wide spectrum of prediction. How is this performance possible despite breast cancer heterogeneity, revealed at the molecular level by numerous and concordant gene profiling studies? We propose that the information contained in our 38-GES relates to a common molecular background, in node-positive and node-negative patients, linked to bad outcome prediction, on the basis of all of our evaluation results, and especially to those obtained by using node-negative cohorts (e.g., Wang's cohort). Furthermore, the biological relevance of the majority of the selected genes strengthens this hypothesis. Disparity of prediction of our 38-GES could let us conclude that this GES is dominated by bad prognosis-associated genes which are also predictive of general chemotherapy sensitivity [29].

**Acknowledgments** This study was supported by SANOFI-AVENTIS-France, PFIZER-France and Cancéropôle Grand Ouest. A part of the tissues used in this work was provided by Institut Régional du Cancer Nantes-Atlantique tumor bank, funded by the Institut National du Cancer and the Cancéropôle Grand Ouest. We thank M. Martin, E. Ollivier, N. Roi and E. Beguet for technical assistance.

## References

- Perou CM, Sorlie T, Eisen MB et al (2000) Molecular portraits of human breast tumours. *Nature* 406:747–752. doi:10.1038/35021093
- Sorlie T, Perou CM, Tibshirani R et al (2001) Gene expression patterns of breast carcinomas distinguish tumor subclasses with clinical implications. *Proc Natl Acad Sci USA* 98:10869–10874. doi:10.1073/pnas.191367098
- Van't Veer LJ, Dai H, van de Vijver MJ et al (2002) Gene expression profiling predicts clinical outcome of breast cancer. *Nature* 415:530–536. doi:10.1038/415530a
- Van de Vijver MJ, He YD, van't Veer LJ et al (2002) A gene-expression signature as a predictor of survival in breast cancer. *N Engl J Med* 347:1999–2009. doi:10.1056/NEJMoa021967
- Wang Y, Klijn JGM, Zhang Y et al (2005) Gene-expression profiles to predict distant metastasis of lymph-node-negative primary breast cancer. *Lancet* 365:671–679
- Paik S, Shak S, Tang G et al (2004) A multigene assay to predict recurrence of tamoxifen-treated, node-negative breast cancer. *N Engl J Med* 351:2817–2826. doi:10.1056/NEJMoa041588
- Ma XJ, Wang Z, Ryan PD et al (2004) A two-gene expression ratio predicts clinical outcome in breast cancer patients treated with tamoxifen. *Cancer Cell* 5:607–616. doi:10.1016/j.ccr.2004.05.015
- Roché H, Fumoleau P, Spielmann M et al (2006) Sequential adjuvant epirubicin-based and docetaxel chemotherapy for node-positive breast cancer patients: the PACS01 trial. *J Clin Oncol* 24:5664–5671. doi:10.1200/JCO.2006.07.3916
- Campone M, Campion L, Roché H et al (2008) Prediction of metastatic relapse in node-positive breast cancer: establishment

- of a clinicogenomic model after FEC100 adjuvant regimen. *Breast Cancer Res Treat* 109:491–501. doi:[10.1007/s10549-007-9673-x](https://doi.org/10.1007/s10549-007-9673-x)
10. Sotiriou C, Neo SY, McShane LM et al (2003) Breast cancer classification and prognosis based on gene expression profiles from a population-based study. *Proc Natl Acad Sci USA* 100:10393–10398. doi:[10.1073/pnas.1732912100](https://doi.org/10.1073/pnas.1732912100)
  11. Sotiriou C, Wirapati P, Loi S et al (2006) Gene expression profiling in breast cancer: understanding the molecular basis of histologic grade to improve prognosis. *J Natl Cancer Inst* 98:262–272
  12. Pawitan Y, Bjohle J, Amler L et al (2005) Gene expression profiling spares early breast cancer patients from adjuvant therapy: derived and validated in two population-based cohorts. *Breast Cancer Res* 7:R953–R964. doi:[10.1186/bcr1325](https://doi.org/10.1186/bcr1325)
  13. Ivshina AV, George J, Senko O et al (2006) Genetic reclassification of histologic grade delineates new clinical subtypes of breast cancer. *Cancer Res* 66:10292–10301. doi:[10.1158/0008-5472.CAN-05-4414](https://doi.org/10.1158/0008-5472.CAN-05-4414)
  14. Chirgwin JM, Przybyla AE, MacDonald RJ et al (1979) Isolation of biologically active ribonucleic acid from sources enriched in ribonuclease. *Biochemistry* 18:5294–5299. doi:[10.1021/bi00591a005](https://doi.org/10.1021/bi00591a005)
  15. Magrangeas F, Nasser V, Avet-Loiseau H et al (2003) Gene expression profiling of multiple myeloma reveals molecular portraits in relation to the pathogenesis of the disease. *Blood* 101:4998–5006. doi:[10.1182/blood-2002-11-3385](https://doi.org/10.1182/blood-2002-11-3385)
  16. Decaux O, Lodé L, Magrangeas F et al (2008) Prediction of survival in multiple myeloma based on gene-expression profiles revealed cell cycle and chromosomal instability signatures in high-risk patients and hyperdiploid signatures in low-risk patients. *J Clin Oncol* 26:4798–4805. doi:[10.1200/JCO.2007.13.8545](https://doi.org/10.1200/JCO.2007.13.8545)
  17. Heagerty PJ, Lumley T, Pepe MS (2000) Time-dependent ROC curves for recurrence and 15-survival data and a diagnostic marker. *Biometrics* 56:337–344. doi:[10.1111/j.0006-341X.2000.00337.x](https://doi.org/10.1111/j.0006-341X.2000.00337.x)
  18. Dupuy A, Simon RM (2007) Critical review of published microarray studies for cancer outcome and guidelines on statistical analysis and reporting. *J Natl Cancer Inst* 99:147–157. doi:[10.1093/jnci/djk018](https://doi.org/10.1093/jnci/djk018)
  19. Simon R (2005) Roadmap for developing and validating therapeutically relevant genomic classifiers. *J Clin Oncol* 23:7332–7341. doi:[10.1200/JCO.2005.02.8712](https://doi.org/10.1200/JCO.2005.02.8712)
  20. Tinker AV, Boussioutas A, Bowtell DDL (2006) The challenges of gene expression microarrays for the study of human cancer. *Cancer Cell* 9:333–339. doi:[10.1016/j.ccr.2006.05.001](https://doi.org/10.1016/j.ccr.2006.05.001)
  21. Sparano JA, Paik S (2008) Development of the 21-gene assay and its application in clinical practice and clinical trials. *J Clin Oncol* 26:721–728. doi:[10.1200/JCO.2007.15.1068](https://doi.org/10.1200/JCO.2007.15.1068)
  22. Habel LA, Shak S, Jacobs MK et al (2006) A population-based study of tumor gene expression and risk of breast cancer death among lymph node-negative patients. *Breast Cancer Res* 8:R25. doi:[10.1186/bcr1412](https://doi.org/10.1186/bcr1412)
  23. Gianni L, Zambetti M, Clark K et al (2005) Gene expression profiles in paraffin-embedded core biopsy tissue predict response to chemotherapy in women with locally advanced breast cancer. *J Clin Oncol* 23:7265–7277. doi:[10.1200/JCO.2005.02.0818](https://doi.org/10.1200/JCO.2005.02.0818)
  24. Fisher B, Dignam J, Wolmark N et al (1997) Tamoxifen and chemotherapy for lymph node-negative, estrogen receptor-positive breast cancer. *J Natl Cancer Inst* 89:1673–1682. doi:[10.1093/jnci/89.22.1673](https://doi.org/10.1093/jnci/89.22.1673)
  25. Paik S, Tang G, Shak S et al (2006) Gene expression and benefit of chemotherapy in women with node-negative, Estrogen receptor-positive breast cancer. *J Clin Oncol* 24:3726–3734. doi:[10.1200/JCO.2005.04.7985](https://doi.org/10.1200/JCO.2005.04.7985)
  26. Albain K, Barlow W, Shak S et al (2007) Prognostic and predictive value of the 21-gene recurrence score assay in postmenopausal, node-positive, ER-positive breast cancer (S8814,INT0100). Presented at 30th Annual San Antonio Breast Cancer Symposium, San Antonio
  27. Goldstein LJ, Gray R, Badve S et al (2008) Prognostic utility of the 21-gene assay in hormone receptor positive operable breast cancer compared with Adjuvant! an integrator of clinicopathologic and treatment information. *J Clin Oncol* 26:4063–4071
  28. Van de Vijver MJ (2007) Gene-expression profiling and the future of adjuvant therapy. *Oncologist* 10(Suppl 2):30–34. doi:[10.1634/theoncologist.10-90002-30](https://doi.org/10.1634/theoncologist.10-90002-30)
  29. Pusztai L, Anderson K, Hess KR (2007) Pharmacogenomic predictor discovery in phase II clinical trials for breast cancer. *Clin Cancer Res* 13:6080–6086. doi:[10.1158/1078-0432.CCR-07-0809](https://doi.org/10.1158/1078-0432.CCR-07-0809)