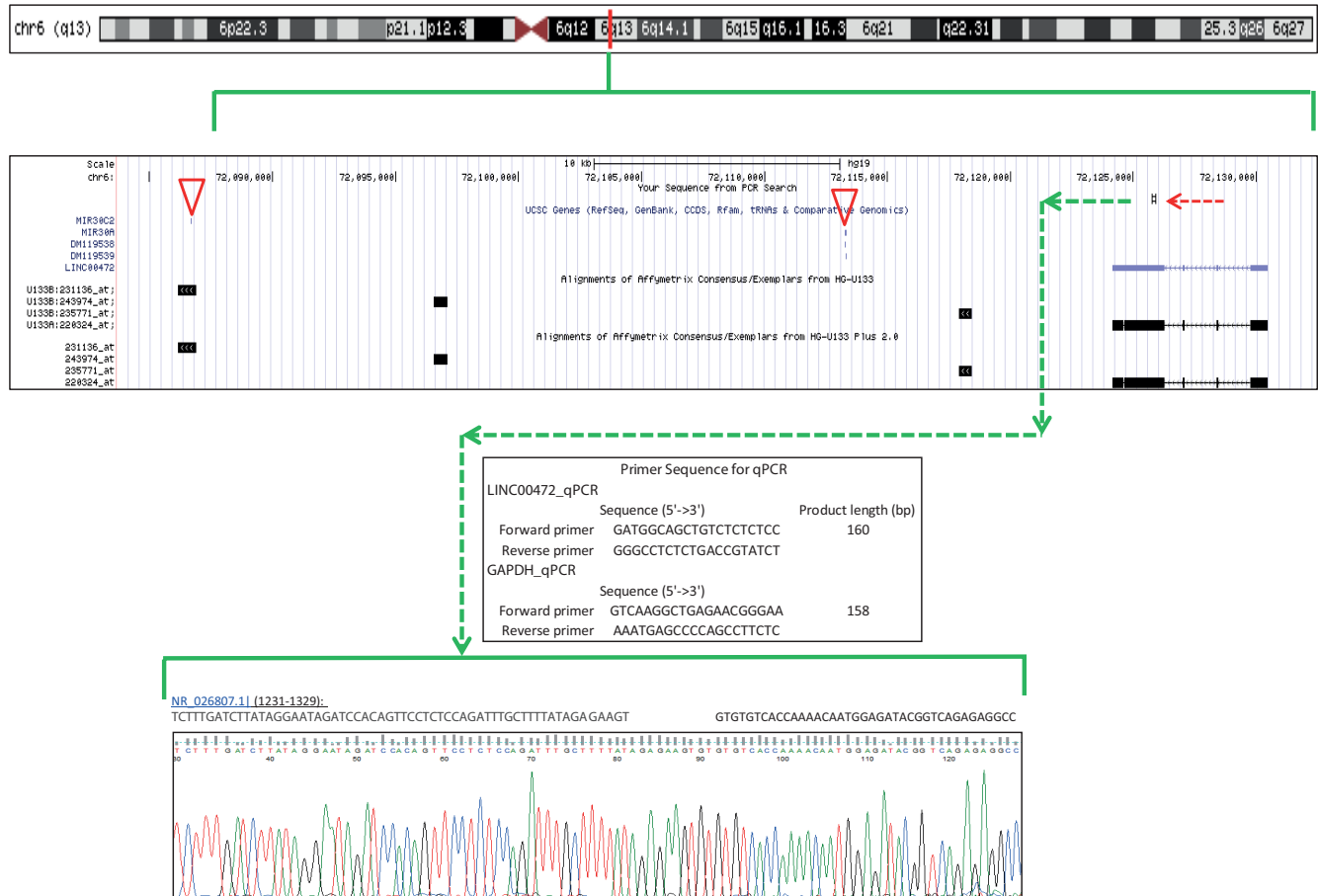


SUPPLEMENTARY FIGURES AND TABLE



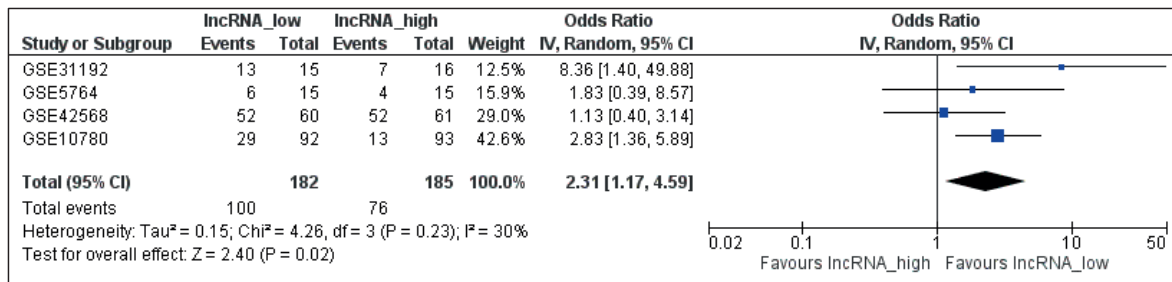
Supplementary Figure S1: *LINC00472* location in chromosome 6 and methods of detection. A schematic diagram shows that *LINC00472* gene is located on chromosome 6q13, complementary strand (red bar). This region is amplified in the following part of the figure. Two microRNA (pointed by red triangles), miR30A and miR30C2, are in the downstream of *LINC00472*. The *LINC00472* transcript (NR_026807.1) includes 4 exons. The red arrow denotes the targeted area of the primers specifically designed for RT-qPCR assay. The black blocks show the targeted area of the probes designed in Affymetrix array HG-U133 and HG-U133 plus 2. The primers' sequences for *LINC00472* and control GAPDH are listed underneath. The RT-qPCR products had been purified and sequenced. The blast result of the PCR product is at the bottom of the figure, showing its match to *LINC00472* transcript NR_026807.1.

A

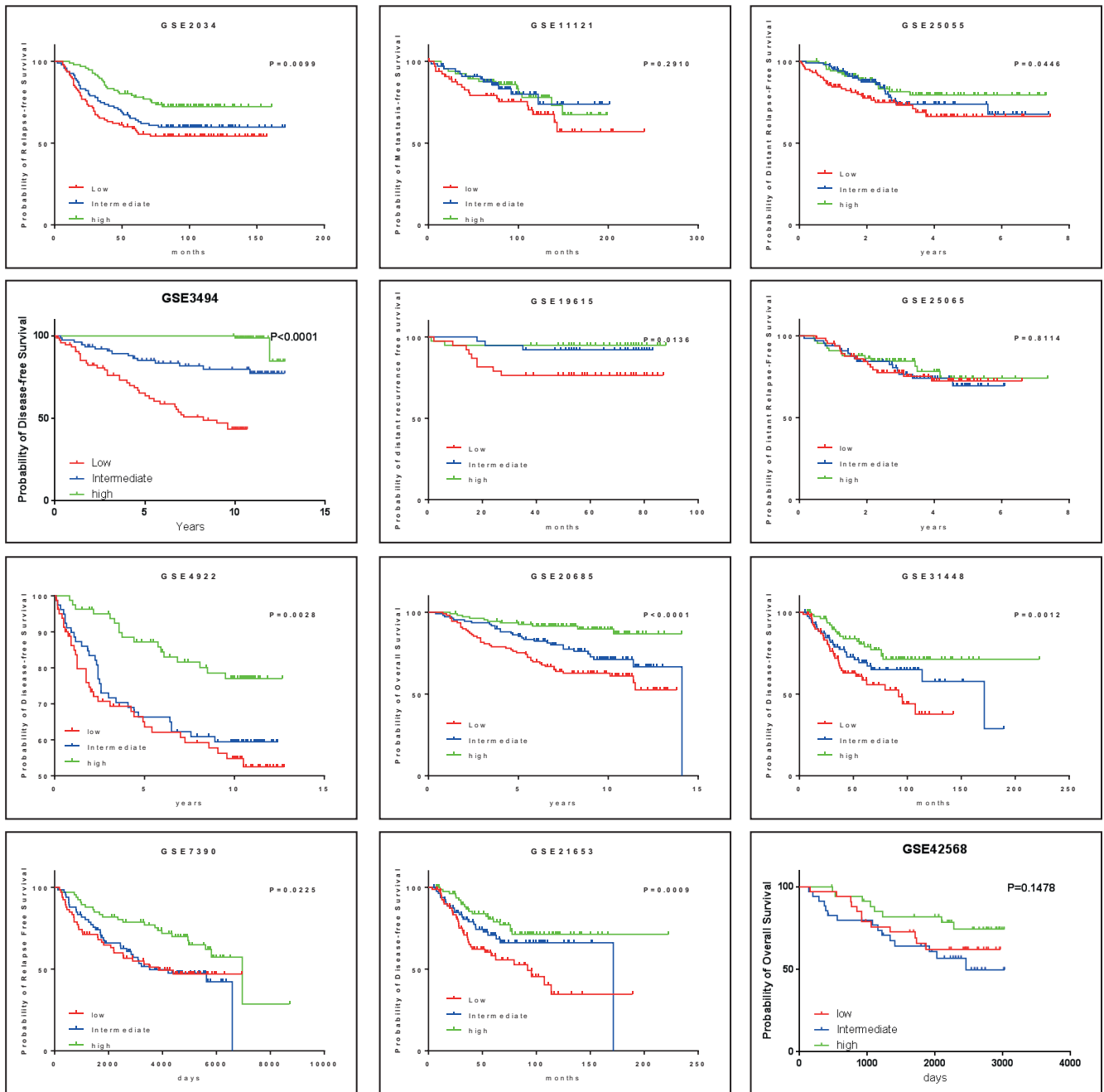
GEO Series ^(Ref)	Country	Normal	Breast Cancer	Platform*	Analysis Method
GSE5764 ¹	Czech Republic	20	10	GPL570	RMA log 2 expression levels
GSE10780 ²	USA	143	42	GPL570	RMA log 2 expression levels
GSE31192 ³	USA	13	20	GPL570	RMA log 2 expression levels
GSE42568 ⁴	Ireland	17	104	GPL570	Log2 GC-RMA signal intensity
Total		193	176		

1. Turashvili G, Bouchal J, Baumforth K, et al. Novel markers for differentiation of lobular and ductal invasive breast carcinomas by laser microdissection and microarray analysis. BMC cancer 2007;7:55.
 2. Chen DT, Nasir A, Culhane A, et al. Proliferative genes dominate malignancy-risk gene signature in histologically-normal breast tissue. Breast cancer research and treatment 2010;119:335-46.
 3. & DMEHJKJOB, & A-CTVFBPS, Horwitz BMJKB. Genomic Signatures of Pregnancy-Associated Breast Cancer Epithelia and Stroma and their Regulation by Estrogens and Progesterone. HORM CANC 2013;4:140-53.
 4. Clarke C, Madden SF, Doolan P, et al. Correlating transcriptional networks to breast cancer survival: a large-scale coexpression analysis. Carcinogenesis 2013;34:2300-8.
 * GPL570 - Affymetrix Human Genome U133 Plus 2.0 Array

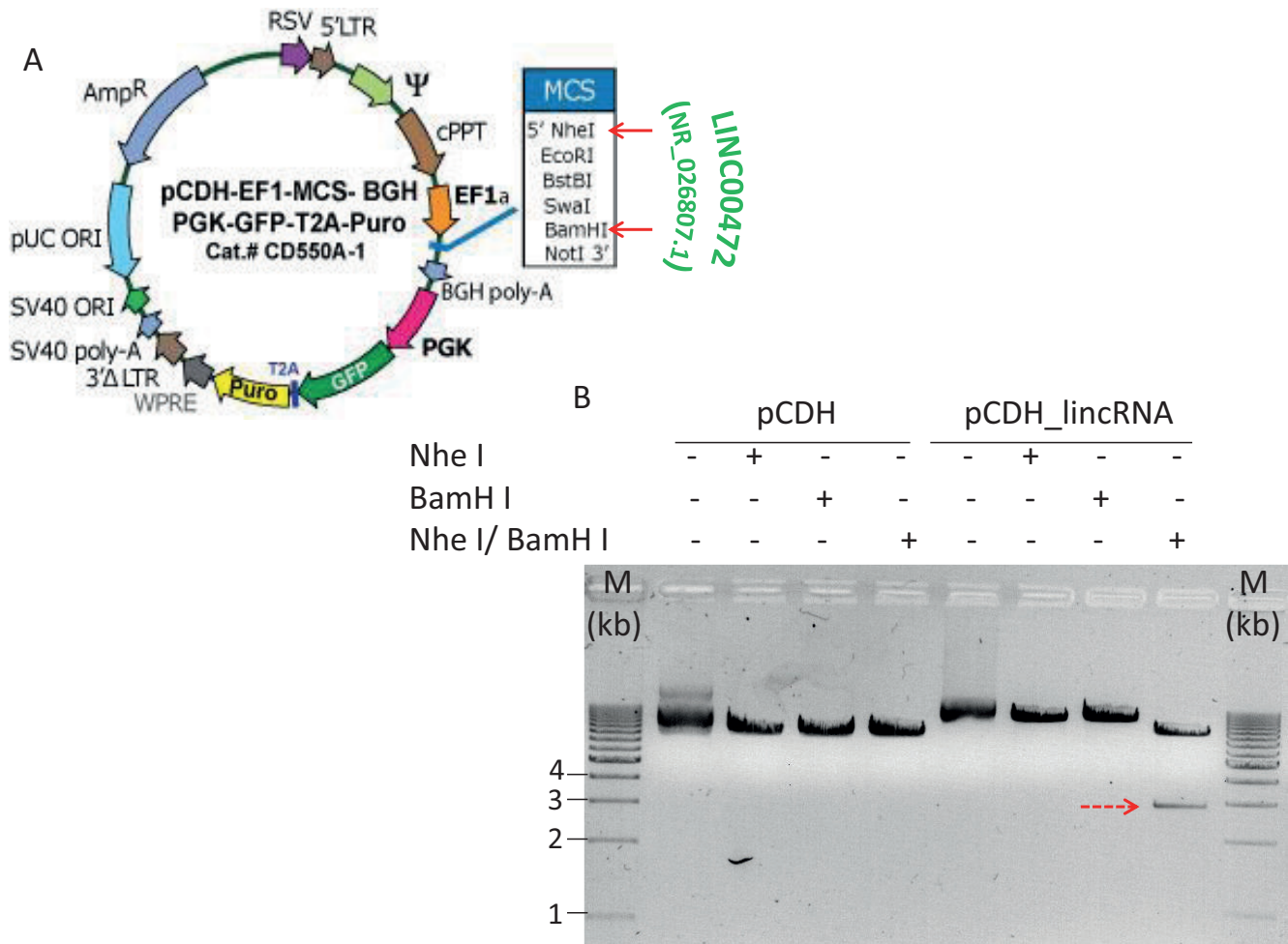
B



Supplementary Figure S2: *LINC00472* expression in normal breast and tumor tissues. A. Study information on 4 GEO datasets which consist of normal breast and breast cancer samples. B. Meta-analysis of the 4 GEO datasets showed higher levels of *LINC00472* expression in normal than in tumor tissues.

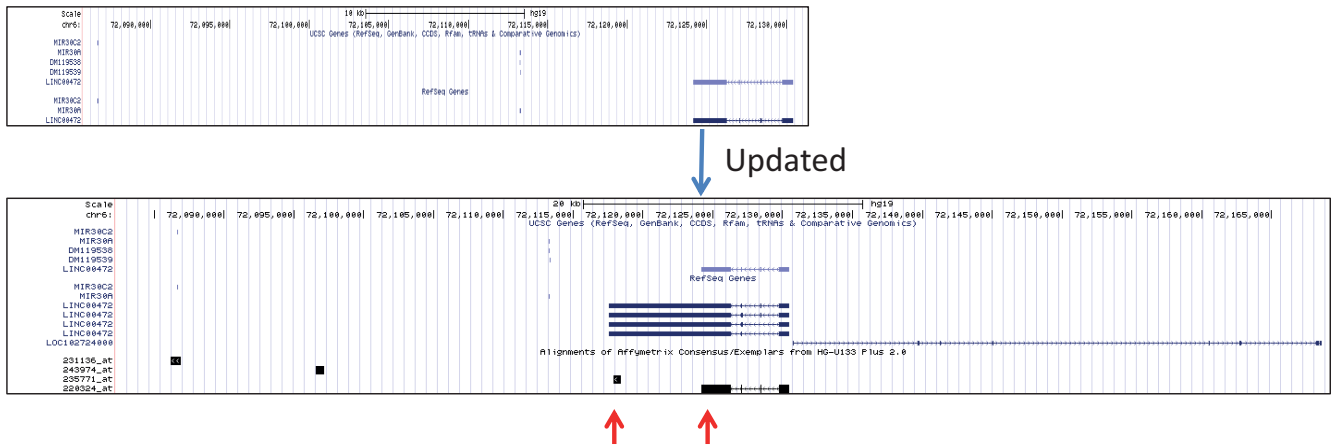


Supplementary Figure S3: Kaplan-meier survival curves by low, medium and high *LINC00472* expression in 12 studies from the GEO database.

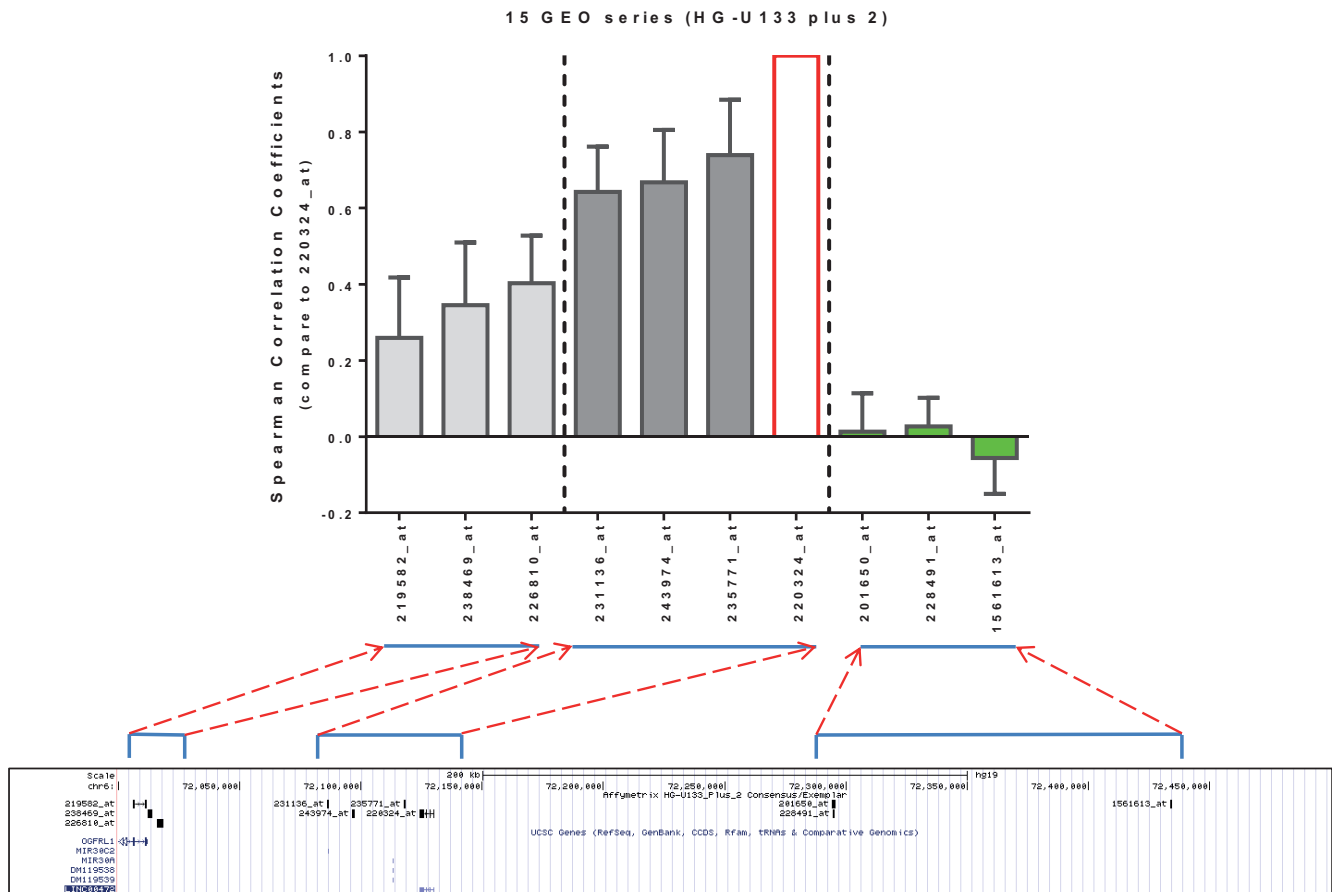


Supplementary Figure S4: *LINC00472* construct and verification. **A.** A *LINC00472* transcript (2933bp, NR_026807.1) was assembled from several EST clones (EHS1001-207275390, EHS1001-207498495, EHS1001-207533792, EHS1001-207590772, EHS1001-210281579, and EHS1001-211231922 from Thermo Scientific Open Biosystems). Two restriction sites, NheI and BamHI, were incorporated into the 5' and 3' end of the transcript, respectively. The entire sequence was inserted into a lentiviral expression vector, pCDH-EF1-MCS-pA-PGK-copGFP-T2A-Puro (System Biosciences). In this vector, EF1a promoter drives the expression of *LINC00472* transcript and PGK promoter directs a report gene, GFP, serving as a transfection control. **B.** Restriction enzyme digestion was performed to the control vector, pCDH and the vector with *LINC00472* transcript, pCDH_ *LINC00472*. Agarose gel electrophoresis showed that a single digestion produced only one DNA fragment both in pCDH and pCDH_ *LINC00472*, and a double digestion on pCDH_ *LINC00472* produced an additional DNA fragment at about a 3 kb position (red arrow), indicating the presence of the inserted transcript, which didn't show in the double digestion of pCDH.

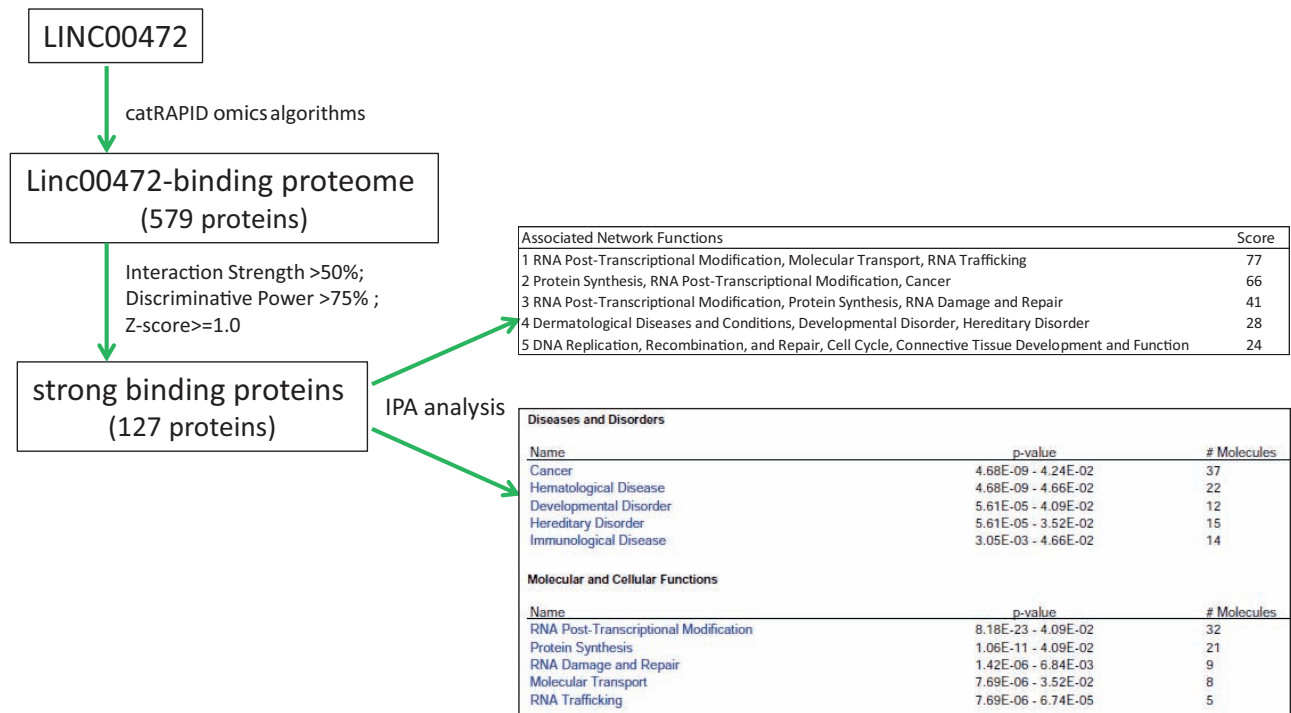
UCSC Genome Browser on Human Feb. 2009 (GRCh37/hg19) Assembly



Supplementary Figure S5: Updated *LINC00472* transcripts. Recently *LINC00472* transcripts have been updated to include NR_121612.1, NR_026807.2, NR_121613.1 and NR_121614.1. The black blocks show the targeted areas of the probes designed in Affymetrix array HG-U133 plus 2, and the red arrows indicate the probes from which we used expression data for analysis.

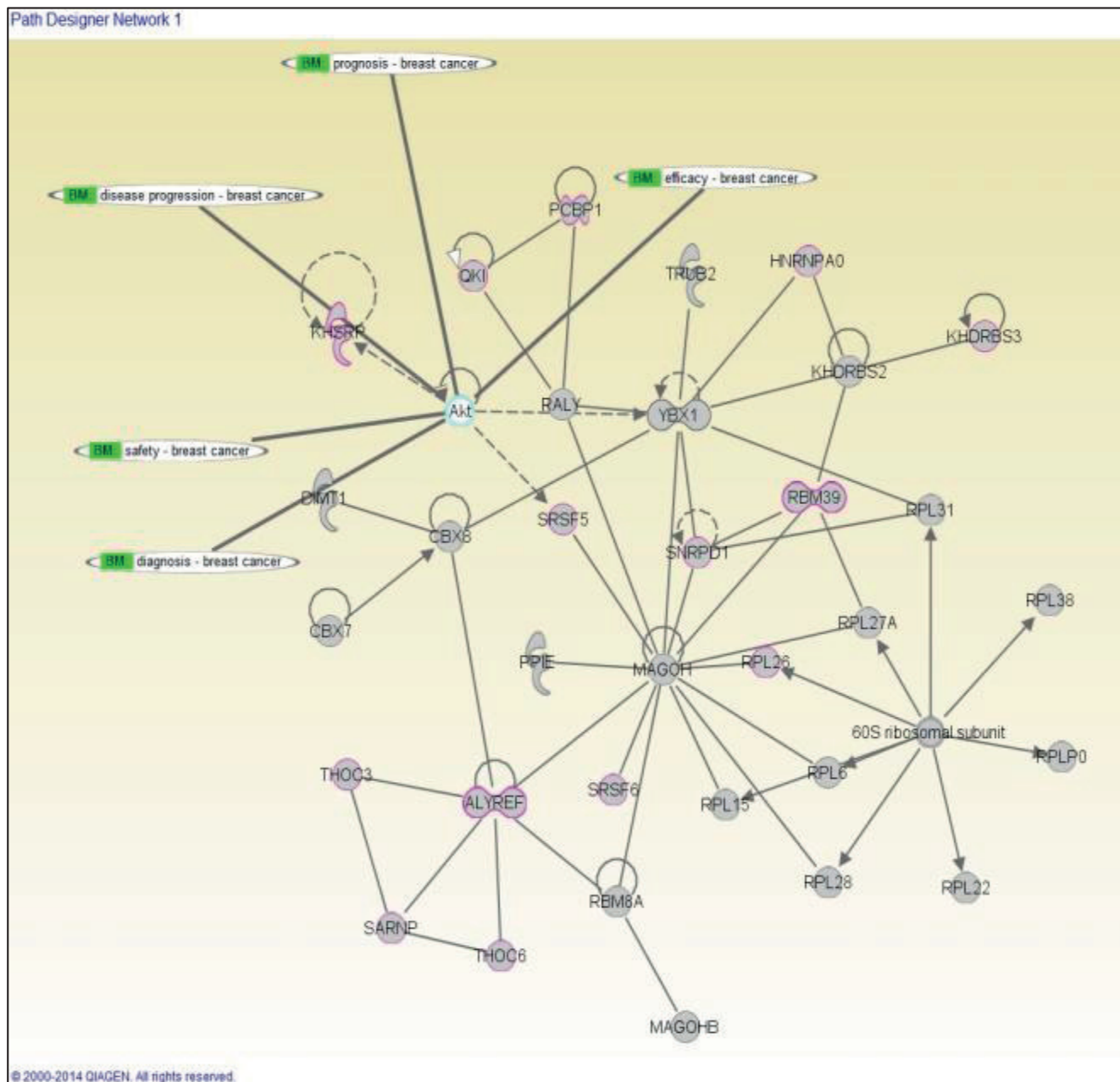


Supplementary Figure S6: Multiple probes in HG-U133 plus 2 array for *LINC00472*. The probe, 235771_at (in dark grey), covers the updated *LINC00472* transcripts as probe 220324_at (red rectangle) does. The expression intensity of these two probes is strongly correlated in 15 GEO datasets. The values in y-axis show the correlation coefficients between expression intensity from different probes and that from probe 220324_at. The grey and green bars represent the probes further away from probe 220324_at.

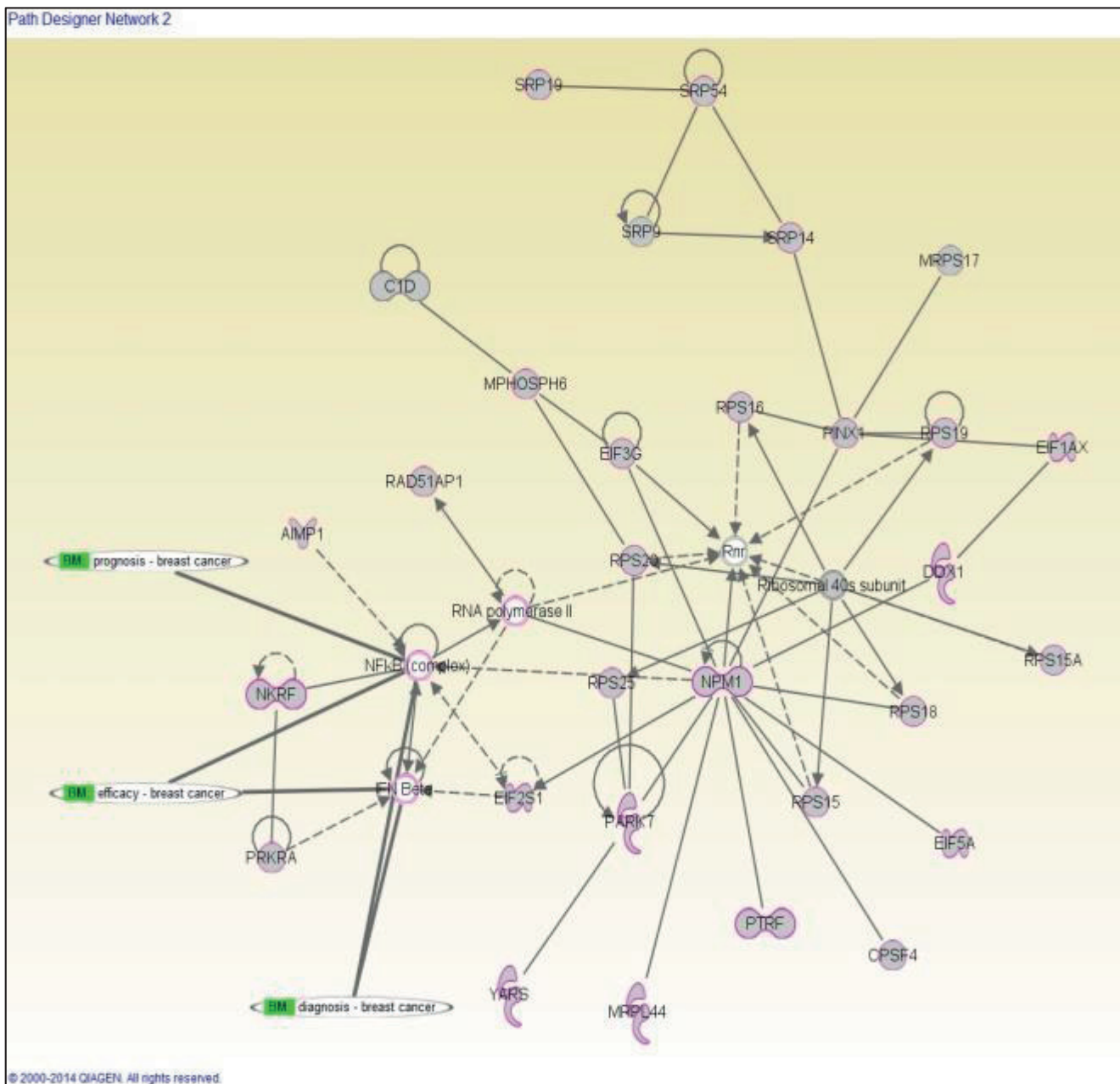


Supplementary Figure S7: Prediction of *LINC00472*-binding proteins. *catRAPID omics* algorithm predicts that *LINC00472* has 579 binding proteins, in which 127 proteins have strong interaction propensities. Ingenuity Core Analysis with these proteins shows the top disease associated with *LINC00472*-binding proteins is cancer, and the top associated network of *LINC00472*-binding proteins is RNA post-transcriptional modification.

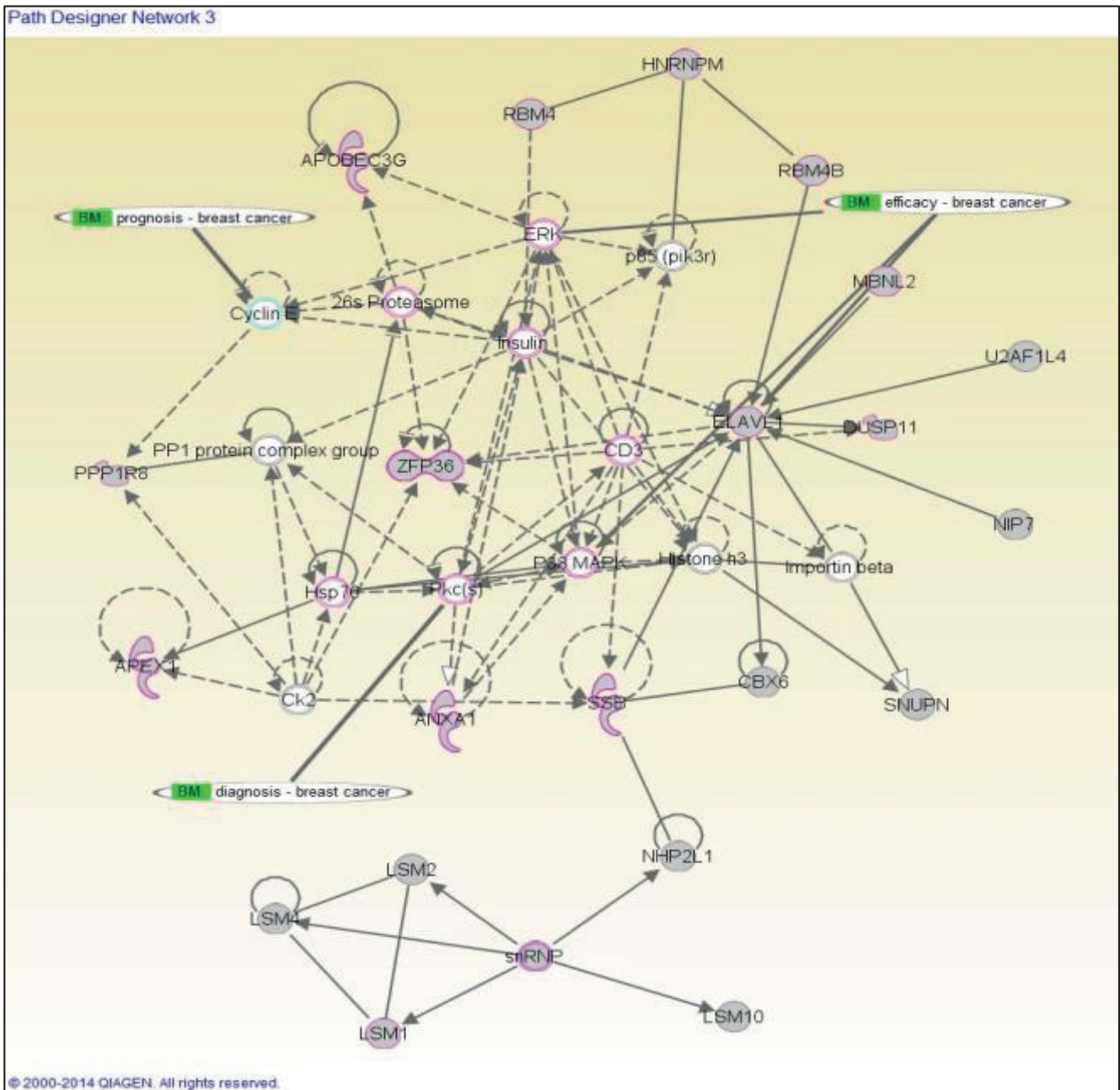
Supplementary Figure S8A–S8C: Top 3 networks associated with *LINC00472*-binding proteins. IPA analysis showed 3 networks (S8A, S8B and S8C) relevant to *LINC00472*. The solid lines denote a robust correlation with partner proteins, and dashed lines indicate statistically significant but less frequent correlations. The grey nodes indicate the *LINC00472*-binding proteins predicted by *catRAPID*; the red color represents the proteins associated with RNA post-transcriptional modifications, Protein Synthesis and Cancer; the un-colored nodes indicate additional proteins of this network that were not spotted by the prediction. The shape of the nodes denote the protein's function: enzymes (diamond); nuclear receptors (rectangle); transcription regulators (oval); cytokines (square); transporter (trapezoid); and others (circle). The breast cancer-related biomarkers found in the networks are marked with BM in green.



(Continued)



(Continued)



Supplementary Table S1. Summary of Publically Available Breast Cancer Microarray Datasets

GEO Series (Ref)	Country	Case Number	Clinical Information						Platform	Analysis Method
			ER	PR	Her2	Grade	Nodal status	Molecular subtype		
GSE2034 ¹	Netherlands	286	Yes						GPL96 ^f	log2 transformed MAS5.0 signal intensity
GSE2603 ²	USA	99	Yes	Yes	Yes				GPL96	log2 transformed MAS5.0 signal intensity
GSE3494 ³	Sweden	251	Yes	Yes		Yes			GPL96	log2 transformed MAS5.0 signal intensity
GSE4922 ⁴	Sweden, Singapore	289	Yes			Yes			GPL96	log2 transformed MAS5.0 signal intensity
GSE5460 ⁵	USA	129	Yes		Yes	Yes			GPL570 ^g	dChip signal intensity
GSE6532 ⁶⁻⁸	UK, Sweden	87		Yes		Yes			GPL570	RMA log 2 expression levels
GSE7390 ⁹	France, UK, Sweden	198	Yes			Yes			GPL96	log2 transformed MAS5.0 signal intensity
GSE11121 ¹⁰	Germany	200				Yes			GPL96	log2 transformed MAS5.0 signal intensity
GSE18864 ¹¹⁻¹³	USA	84	Yes	Yes	Yes	Yes			GPL570	RMA log 2 expression levels
GSE19615 ¹²	USA	115	Yes	Yes	Yes	Yes	Yes		GPL570	dChip signal intensity
GSE20194 ^{14,15}	USA	278	Yes	Yes	Yes	Yes			GPL96	log2 transformed MAS5.0 signal intensity
GSE20271 ¹⁶	USA	178	Yes	Yes	Yes	Yes			GPL96	log10 transformed dChip signal intensity
GSE20685 ¹⁷	Taiwan	327							GPL570	log2 transformed MAS5.0 signal intensity
GSE21653 ^{18,19}	France	266	Yes	Yes		Yes	Yes		GPL570	RMA log 2 expression levels
GSE23177 ²⁰	Belgium	116					Yes		GPL570	log2 transformed MAS5.0 signal intensity

(Continued)

GEO Series (Ref)	Country	Case Number	Clinical Information							Platform	Analysis Method
			ER	PR	Her2	Grade	Nodal status	Molecular subtype	Survival outcomes		
GSE23720 ^{18,19}	France	197	Yes	Yes						GPL570	RMA log 2 expression levels
GSE23988 ²¹	USA	61	Yes			Yes		Yes		GPL96	log2 transformed dChip signal intensity
GSE24185 ²²	USA	103	Yes	Yes	Yes			Yes		GPL96	dChip signal intensity
GSE25055 ²³	USA	310	Yes	Yes	Yes	Yes	Yes	Yes		GPL96	log2 transformed MAS5.0 signal intensity
GSE25065 ²³	USA	198	Yes	Yes	Yes	Yes	Yes	Yes		GPL96	log2 transformed MAS5.0 signal intensity
GSE31448 ^{18,19}	France	353	Yes	Yes		Yes	Yes			GPL570	RMA log 2 expression levels
GSE42568 ²⁴	Ireland	104	Yes			Yes	Yes	Yes		GPL570	Log2 GC-RMA signal intensity
GSE47109 ²⁵	USA	246	Yes				Yes			GPL570	log2 transformed MAS5.0 signal intensity
GSE48390 ²⁶	Taiwan	81	Yes		Yes					GPL570	RMA log 2 expression levels
Turin_study	Italy	348	Yes	Yes		Yes	Yes	Yes		Q-PCR	Expression index
Total:		4904									

^aDMFS - Distance Metastasis Free Survival

^bLMFS - Lung Metastasis Free Survival

^cBMFS - Bone Metastasis Free Survival

^dDFS - Disease Free Survival

^eOS - Overall Survival

^fGPL96 - Affymetrix Human Genome U133A Array (HG-U133A)

^gGPL570 - Affymetrix Human Genome U133 Plus 2.0 Array (HG-U133 Plus 2)

REFERENCES IN SUPPLEMENTARY TABLE S1

- Wang Y, Klijn JG, Zhang Y, et al. Gene-expression profiles to predict distant metastasis of lymph-node-negative primary breast cancer. *Lancet*. 2005; 365:671–9.
- Minn AJ, Gupta GP, Siegel PM, et al. Genes that mediate breast cancer metastasis to lung. *Nature*. 2005; 436:518–24.
- Miller LD, Smeds J, George J, et al. An expression signature for p53 status in human breast cancer predicts mutation status, transcriptional effects, and patient survival. *Proceedings of the National Academy of Sciences of the United States of America*. 2005; 102:13550–5.
- Ivshina AV, George J, Senko O, et al. Genetic reclassification of histologic grade delineates new clinical subtypes of breast cancer. *Cancer research*. 2006; 66:10292–301.
- Lu X, Lu X, Wang ZC, Iglehart JD, Zhang X, Richardson AL. Predicting features of breast cancer with gene expression patterns. *Breast cancer research and treatment*. 2008; 108:191–201.
- Loi S, Haibe-Kains B, Desmedt C, et al. Definition of clinically distinct molecular subtypes in estrogen receptor-positive breast carcinomas through genomic grade. *Journal of clinical oncology : official journal of the American Society of Clinical Oncology*. 2007; 25:1239–46.
- Loi S, Haibe-Kains B, Majjaj S, et al. PIK3CA mutations associated with gene signature of low mTORC1 signaling and better outcomes in estrogen receptor-positive breast cancer. *Proceedings of the National Academy of Sciences of the United States of America*. 2010; 107:10208–13.
- Loi S, Haibe-Kains B, Desmedt C, et al. Predicting prognosis using molecular profiling in estrogen receptor-positive breast cancer treated with tamoxifen. *BMC genomics*. 2008; 9:239.
- Desmedt C, Piette F, Loi S, et al. Strong time dependence of the 76-gene prognostic signature for node-negative breast cancer patients in the TRANSBIG multicenter independent validation series. *Clinical cancer research: an official journal of the American Association for Cancer Research*. 2007; 13:3207–14.
- Schmidt M, Bohm D, von Torne C, et al. The humoral immune system has a key prognostic impact in node-negative breast cancer. *Cancer research*. 2008; 68:5405–13.
- Silver DP, Richardson AL, Eklund AC, et al. Efficacy of neoadjuvant Cisplatin in triple-negative breast cancer. *Journal of clinical oncology: official journal of the American Society of Clinical Oncology*. 2010; 28:1145–53.
- Li Y, Zou L, Li Q, et al. Amplification of LAPT4B and YWHAZ contributes to chemotherapy resistance and recurrence of breast cancer. *Nature medicine*. 2010; 16:214–8.
- Juul N, Szallasi Z, Eklund AC, et al. Assessment of an RNA interference screen-derived mitotic and ceramide pathway metagene as a predictor of response to neoadjuvant paclitaxel for primary triple-negative breast cancer: a retrospective analysis of five clinical trials. *The lancet oncology*. 2010; 11:358–65.
- Popovici V, Chen W, Gallas BG, et al. Effect of training-sample size and classification difficulty on the accuracy of genomic predictors. *Breast cancer research: BCR*. 2010; 12:R5.
- Shi L, Campbell G, Jones WD, et al. The MicroArray Quality Control (MAQC)-II study of common practices for the development and validation of microarray-based predictive models. *Nature biotechnology*. 2010; 28:827–38.
- Tabchy A, Valero V, Vidaurre T, et al. Evaluation of a 30-gene paclitaxel, fluorouracil, doxorubicin, and cyclophosphamide chemotherapy response predictor in a multicenter randomized trial in breast cancer. *Clinical cancer research: an official journal of the American Association for Cancer Research*. 2010; 16:5351–61.
- Kao KJ, Chang KM, Hsu HC, Huang AT. Correlation of microarray-based breast cancer molecular subtypes and clinical outcomes: implications for treatment optimization. *BMC cancer*. 2011; 11:143.
- Sabatier R, Finetti P, Adelaide J, et al. Down-regulation of ECRG4, a candidate tumor suppressor gene, in human breast cancer. *PloS one*. 2011; 6:e27656.
- Sabatier R, Finetti P, Cervera N, et al. A gene expression signature identifies two prognostic subgroups of basal breast cancer. *Breast cancer research and treatment*. 2011; 126:407–20.
- Smeets A, Daemen A, Vanden Bempt I, et al. Prediction of lymph node involvement in breast cancer from primary tumor tissue using gene expression profiling and miRNAs. *Breast cancer research and treatment*. 2011; 129:767–76.
- Iwamoto T, Bianchini G, Booser D, et al. Gene pathways associated with prognosis and chemotherapy sensitivity in molecular subtypes of breast cancer. *Journal of the National Cancer Institute*. 2011; 103:264–72.
- Creighton CJ, Sada YH, Zhang Y, et al. A gene transcription signature of obesity in breast cancer. *Breast cancer research and treatment*. 2012; 132:993–1000.
- Hatzis C, Pusztai L, Valero V, et al. A genomic predictor of response and survival following taxane-anthracycline chemotherapy for invasive breast cancer. *JAMA: the journal of the American Medical Association*. 2011; 305:1873–81.
- Clarke C, Madden SF, Doolan P, et al. Correlating transcriptional networks to breast cancer survival: a large-scale coexpression analysis. *Carcinogenesis*. 2013; 34:2300–8.
- D'Alfonso TM, van Laar RK, Vahdat LT, et al. BreastPRS is a gene expression assay that stratifies intermediate-risk Oncotype DX patients into high- or low-risk for disease recurrence. *Breast cancer research and treatment*. 2013; 139:705–15.
- Huang CC, Tu SH, Lien HH, et al. Concurrent gene signatures for han chinese breast cancers. *PloS one*. 2013; 8:e76421.