

Molecular Profiles of Matched Primary and Metastatic Tumor Samples Support a Linear Evolutionary Model of Breast Cancer

Runpu Chen¹, Steve Goodison², and Yijun Sun^{1,3,4}



ABSTRACT

The interpretation of accumulating genomic data with respect to tumor evolution and cancer progression requires integrated models. We developed a computational approach that enables the construction of disease progression models using static sample data. Application to breast cancer data revealed a linear, branching evolutionary model with two distinct trajectories for malignant progression. Here, we used the progression model as a foundation to investigate the relationships between matched primary and metastasis breast tumor samples. Mapping paired data onto the model confirmed that molecular breast cancer subtypes can shift during progression and supported directional

tumor evolution through luminal subtypes to increasingly malignant states. Cancer progression modeling through the analysis of available static samples represents a promising breakthrough. Further refinement of a roadmap of breast cancer progression will facilitate the development of improved cancer diagnostics, prognostics, and targeted therapeutics.

Significance: Analysis of matched primary and metastatic tumor samples supports a unidirectional, linear cancer evolution process and sheds light on longstanding issues regarding the origins of molecular subtypes and their progression relationships.

Introduction

Human cancer is a dynamic disease that develops over an extended period of time through the accumulation of genetic alterations. Once initiated, the advance of the disease to malignancy can be viewed as a Darwinian, multistep evolutionary process at the cellular level (1). Delineating the dynamic disease process and identifying pivotal molecular events that drive stepwise disease progression would significantly advance our understanding of tumorigenesis and provide a foundation for the development of improved cancer diagnostics, prognostics, and targeted therapeutics. Traditionally, system dynamics is approached through time-series studies achieved by repeated sampling of the same cohort of subjects across an entire biological process. However, due to the need for timely surgical intervention upon diagnosis, it is not ethically feasible to collect time-series data to study human cancer. Consequently, while the concept of cancer evolution has been widely accepted (1, 2), the biological process of how cancer progresses to a malignant, life-threatening disease is still not well-understood. The lack of time-series data has been recognized by the field as the central problem in studying cancer progression (3).

With the rapid development of sequencing technology, many thousands of excised tumor tissue specimens are being collected in large-scale

cancer studies (4–6). This provides us with a unique opportunity to develop a novel analytical strategy to use static data, instead of time-course data, to study disease dynamics. The strategy is based on the rationale that each excised tissue sample provides a snapshot of the disease process, and if the number of samples is sufficiently large, the genetic footprints of individual samples populate progression trajectories, enabling us to recover disease dynamics by using computational approaches. We developed a comprehensive bioinformatics pipeline (7) and applied it to the gene expression data from over 3,100 breast tumor samples available from The Cancer Genome Atlas and METABRIC consortiums (4, 5). Our analysis demonstrated that it is indeed possible to use static sample data to study disease dynamics and led to one of the first working models of breast cancer progression that covers the entire disease process. The progression pattern was confirmed by analysis of a series of smaller independent breast cancer datasets and by aligning established clinical and molecular traits with the model (7). To further validate the model, here we proposed a novel strategy to investigate the progression relationships of matched primary and metastasis (P/M) tumor samples from patients with breast cancer. Our analysis suggested that while breast cancer is a genetically and clinically heterogeneous disease, tumor samples are distributed on a low-dimensional manifold, that disease subtypes are not hardwired and can shift within the same individual, and that the shift is unidirectional along a continuum of disease state toward malignancy. This study shed light on some longstanding issues regarding the origins of molecular subtypes and their possible progression relationships.

Materials and Methods

The P/M dataset was downloaded from Gene Expression Omnibus (accession number: GSE92977), which contains the expression levels of 105 breast cancer-related genes and 5 house-keeping genes from 246 matched tumor samples collected from 123 patients with breast cancer (8). To overcome the issue of missing data, we first filtered out 6 genes and 13 samples containing >20% missing values, and then performed missing data imputation on the remaining samples. Specifically, for each sample, we found its 10 nearest neighbors and

¹Department of Computer Science and Engineering, The State University of New York, Buffalo, New York. ²Department of Health Sciences Research, Mayo Clinic, Jacksonville, Florida. ³Department of Microbiology and Immunology, The State University of New York, Buffalo, New York. ⁴Department of Biostatistics, The State University of New York, Buffalo, New York.

Note: Supplementary data for this article are available at Cancer Research Online (<http://cancerres.aacrjournals.org/>).

Corresponding Authors: Yijun Sun, University at Buffalo, 3435 Main Street, Buffalo, NY 14214. Phone: 716-881-1374; E-mail: yijunsun@buffalo.edu; and Steve Goodison, Mayo Clinic, Jacksonville, FL 32224; E-mail: goodison.steven@mayo.edu

Cancer Res 2020;80:170–4

doi: 10.1158/0008-5472.CAN-19-2296

©2019 American Association for Cancer Research.

replaced the missing data in each gene in the sample by the average of the observed values of the gene in the identified nearest neighbors. The P/M dataset also contains a considerable number of outlier samples, which would complicate downstream analysis. Intuitively, an outlier sample should be dissimilar to its nearest neighbors. We exploited this intuition by calculating the average distance between each sample and its 10 nearest neighbors and removing 23 samples with top 10% average distances. By using the PAM50 classifier, we stratified the samples into the five intrinsic molecular subtypes (luminal A, luminal B, HER2⁺, basal, and normal-like). With the subtype information, we performed a one-way ANOVA analysis on individual genes and removed 12 genes that contain little information in discriminating the five molecular subtypes ($P > 0.1$). After data preprocessing, 87 genes and 210 samples including 92 pairs were retained for the further analysis.

We investigated the progression relationships between the matched tumor pairs by mapping the P/M data onto a progression model of breast cancer that we constructed using the METABRIC data (referred to as the METABRIC model hereafter; ref. 7). Because the data sources are not entirely compatible, we first mapped the 87 genes in the P/M cohort back to the 25,160 genes in the METABRIC data and identified 85 genes present in both datasets, and then used ComBat (9) to remove the platform-induced data deviations. The METABRIC model was constructed using 359 genes (7), while the P/M cohort contains only the expression data of 85 genes. To address the issue, we designed a novel computational strategy that projects the P/M data onto a high-dimensional space spanned by the 359 genes used to build the METABRIC model. Specifically, we assumed that the projection relationships between the 85 and 359 genes were shared by both cohorts, learned a projection function by using the METABRIC data, and applied the function to the P/M data. See Supplementary Data for a detailed mathematical derivation. After we obtained the high-dimensional prediction of the P/M data, we projected each sample onto the progression paths identified in the METABRIC model. Here, the projection of a sample is defined as a point on a progression path that is the closest to the sample. By using the mean of the normal samples to represent the origin of cancer progression, we compared the progression distances of metastasis tumors with those of their matched primary tumors. Because the P/M cohort measured only the expression levels of 85 genes and contained considerable measurement errors, many pairs had small differences in their progression distances, likely due to random variations. We devised a way to identify tumor pairs that underwent significant disease progression (Supplementary Data).

Results

Constructing a progression model of breast cancer

As a foundation for the investigation of the relationships between matched tumor pairs, we applied our computational pipeline developed in (7) to the METABRIC data and reconstructed a progression model of breast cancer (see Supplementary Data; ref. 7 for details). The dataset contains the expression profiles of 25,160 genes obtained from 1,989 surgically excised breast tumor samples. Briefly, we first performed supervised learning by using the breast cancer subtypes as class labels and selected 359 disease-related genes. Then, we performed a clustering analysis on the expression measures of the selected genes to detect tumor groups with homogenous gene expression profiles. By using gap statistic (10) and consensus clustering (11), 10 distinct clusters were identified. Finally, we constructed a progression model and represented it as an undirected graph, by using the centroids of the identified clusters as the vertices and connecting them based on the

progression trend inferred from the analysis. **Figure 1A** and **B** show the distribution of the tumor samples in a three-dimensional space supported by the selected genes and a schematic of the constructed model, respectively. Our analysis identified a linear, branching model describing two distinct trajectories to malignancy, either directly to the basal subtype with little deviation, or a stepwise, more indolent path through the luminal subtypes to the HER2⁺ subtype. The two trajectory termini (i.e., HER2⁺ and basal) represent the two most aggressive breast tumor subtypes (12). Significant side branches are also evident for both luminal A/B subtypes. Our results confirmed that molecular subtypes are not hardwired, and genotypes and phenotypes can shift over time (1), and that cancer development follows limited, common progression paths (4).

Mapping of matched primary and metastatic tumor samples

To validate the constructed model, we investigated the interrelationships of the matched primary and metastasis tumor samples with respect to molecular subtype and directional progression. To overcome the issues of missing data and outliers, we performed a series of data preprocessing to remove genes and samples containing excessive missing values, to impute missing data, and to identify outlier data (Materials and Methods section). After data preprocessing, a total of 210 samples including 92 pairs were retained for further analysis. Using the PAM50 classifier (13), we stratified the samples into five intrinsic molecular subtypes, including three normal-like, 79 luminal A, 78 luminal B, 33 HER2⁺, and 17 basal tumors (Supplementary Table S1). It has been reported that normal-like samples could be technical artifacts from high contamination of normal tissue (14), so we removed from further analysis the three pairs that contained normal-like classification. **Figure 2A** presents a Sankey plot showing the subtype changes of matched pairs. We then performed a quantitative analysis of the progression relationships of the primary and metastasis tumors by mapping the data onto the METABRIC model. Because the data sources are not entirely compatible, we first performed batch-effect correction by using ComBat (9) and then applied a novel strategy to map the P/M data onto the 359-dimensional space where the METABRIC model was constructed (Materials and Methods section). **Figure 2B** visualizes the sample distributions of both P/M and METABRIC data, where each sample was color-coded on the basis of its PAM50 label. By using the normal samples as the baseline to represent the origin of cancer progression, we compared the progression distances of each matched pair and reported the results in **Fig. 2C**. After multiple testing correction, a total of 16 pairs were identified with significant positive disease progression (i.e., the progression distance of a primary tumor is significantly smaller than that of the matched metastasis tumor) and one pair with negative disease progression ($FDR \leq 0.1$).

From **Fig. 2A**, we observed that all basal metastatic tumors were derived from basal primary tumors, and all nonbasal metastatic tumors were derived from nonbasal primary tumors. This suggests that basal cancer is a distinct disease entity, supporting the bifurcating structure revealed by the proposed progression model (**Fig. 1**). We also observed that while most paired samples (62%) were of the same molecular subtype, subtypes did shift primarily from a lesser to a more malignant phenotype (i.e., luminal A to B or to HER2⁺ subtype), which aligns well with the METABRIC model. There were several examples of luminal B shifting to luminal A (six pairs) and HER2⁺ shifting to luminal B (one pair), implying that cancer evolution can be bidirectional. However, because the PAM50 system provides only an approximate stratification of breast cancer (15), the P/M cohort measured only the expression levels of 85 genes, and the data contains considerable measurement errors, a qualitative analysis of subtype changes

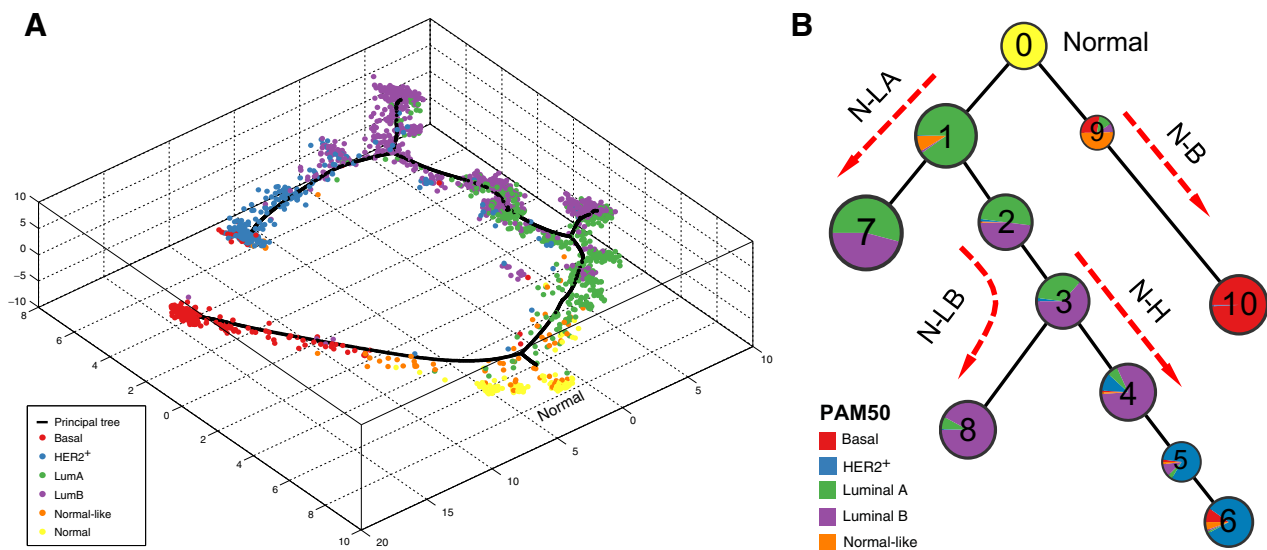


Figure 1.

Breast cancer progression modeling analysis performed on METABRIC gene expression data. **A**, Data visualization analysis provides a general view of sample distribution. The dataset contains 144 normal breast tissue samples, which we used as the baseline to represent the origin of cancer progression. To help with visualization and to put the result into context by referring to previous classification systems, each sample was color-coded on the basis of its PAM50 subtype label. **B**, A progression model of breast cancer. Each node represents an identified cluster, and the pie chart in each node depicts the percentage of the samples in the node belonging to one of the five PAM50 subtypes. The analysis revealed four major progression paths, referred to as N-B (normal to basal), N-H (normal through luminal A/B to HER2⁺), N-LB (normal through luminal A to luminal B side-branch), and N-LA (normal to luminal A side-branch). The model was modified from Sun and colleagues (7). LumA, luminal A; LumB, luminal B.

tells only part of the story. Indeed, as shown in **Fig. 2B**, luminal A and B do not have a clear-cut boundary, as is the case for luminal B and HER2⁺. This boundary overlap of approximate subtypes explains why in some large-scaled benchmark studies it has been observed that existing molecular subtyping methods only achieve moderate concordance, particularly when classifying luminal A and B tumors (16). Through a quantitative analysis of the progression distances revealed by mapping the P/M data to the model (**Fig. 2B** and **C**), we found that all six luminal B primary tumors that changed to luminal A metastases and one HER2⁺ primary tumor that changed to a luminal B metastasis had very small disease progression. Conversely, all 16 pairs with significant positive progression were either luminal A to luminal B (nine pairs), luminal A to HER2⁺ (three), luminal B to HER2⁺ (one), HER2⁺ to HER2⁺ (two), or basal to basal (one). As examples, **Fig. 2B** shows the locations of two pairs (9P/9M and 48P/48M) that underwent evident disease progression. Taken together, these analyses support a unidirectional, linear evolution process for breast cancer through luminal subtypes toward malignancy. Interestingly, one tumor pair (18P/18M, both classified as luminal A) had a significant negative disease progression. Although the overall observed trend here was a downstream P-to-M shift, a number of scenarios can coexist. Tumor cells can escape the primary lesion and survive at early stages of primary tumor development. The evolutionary time index may be different for these foci, especially if the metastasis lesion is dormant for a period, and local selective pressures at primary and secondary sites can differentially impact the evolution of the related lesions.

Discussion

Cancer evolution theory dates back to the 1970s (17), and numerous studies have been conducted that significantly expanded

our understanding of the concept (1, 3, 18). Yet, due to the difficulty in obtaining time-series data, beyond conceptual models, there is currently no established cancer progression model derived from tumor tissue data that covers the entire disease process. We have proposed a novel computational strategy that overcomes the existing sampling restrictions. The application of the approach to large-scale breast cancer genomic datasets identified a linear, bifurcating progression model describing two distinct pathways to cancer malignancy. The interpretation of the model is that the basal subtype is distinct from the luminal subtypes, and that the luminal subtypes can shift during disease progression and may be considered as different stages of the same disease. The mapping onto the progression model of the paired P/M samples further supports the overall model and the concept that the cancer evolutionary process is unidirectional through luminal subtypes toward increasingly malignant states.

Since the proposal of the cancer molecular subtypes (12), fundamental issues regarding whether subtypes are biologically independent entities or have progressive relationships have been under debate (15, 19). One conceptual model proposes a distinct path scenario where each subtype follows a path of initiation and progression independently. The alternative is a linear evolution model, where tumors gradually evolve from normal cells to malignant states through the accumulation of genetic alterations (15). While both models embrace the notion of cancer evolution, the first implies that subtypes are different diseases, while the second suggests that subtypes are different stages of the same disease. Clarifying this issue could have a profound impact on current cancer research because clinical management and research strategies in the two scenarios may be very different. Our analysis supports the second model as a representation of disease progression, but also suggests that basal and luminal/HER2⁺ subtypes are differentially derived from an ancestral somatic or neoplastic cell of origin.

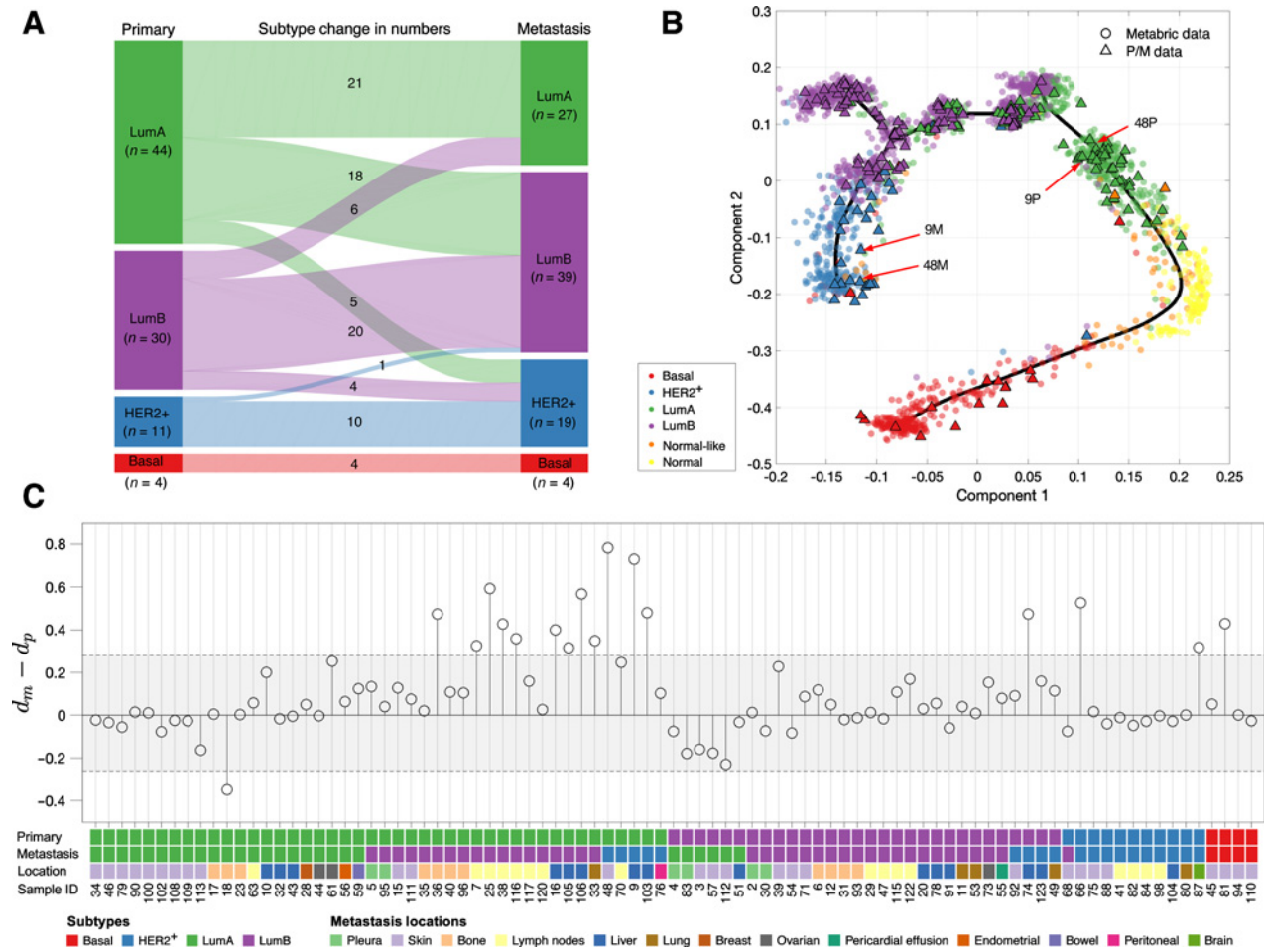


Figure 2. Progression analysis of 89 pairs of matched primary and metastasis tumor samples. **A**, Sankey diagram showing the subtype changes of matched tumor pairs. **B**, Data visualization of the P/M cohort mapped onto the METABRIC model. Two examples (9P/9M and 48P/48M) are shown that underwent evident disease progression from luminal A to the HER2⁺ subtype. **C**, Comparison of progression distances of matched primary and metastasis tumors. A total of 16 and one pairs were identified with significant positive and negative disease progression (i.e., samples outside of the shaded region; FDR ≤ 0.1), respectively. d_m , progression distance of a metastasis tumor; d_p , progression distance of a primary tumor; LumA, luminal A; LumB, luminal B.

The development of cancer progression models can inform a range of research directions. For example, current prognostic tests are of value only in a restricted set of patients, but if we can visualize the entire, ordered progressive process, the identification of specific molecular characteristics associated with a broader spectrum of cancer phenotypes becomes feasible. Assisted by genomic testing, we can envisage the placement of samples from individual cases onto a progression path to guide clinical management and evaluate individualized treatment success. The derivation of annotated progression maps can also guide the design of animal studies to focus on pivotal points of cancer development, which may yield the best return with limited resources. Guided by a working model, future studies using higher resolution sampling (e.g., single-cell sequencing and tissue microdissection) and incorporating multi-omics data can provide refined roadmaps of breast cancer progression. Although in this study we focus mainly on breast cancer, the developed methods for model construction and validation can also be used to study other

distinct (6) or related cancers (20) and other progressive human diseases where the lack of longitudinal data is an unavoidable problem.

Disclosure of Potential Conflicts of Interest

No potential conflicts of interest were disclosed.

Acknowledgments

We thank the editor and the four reviewers for their valuable comments that have helped us significantly improve the quality of the article. We thank Dr. Anthony A. Campagnari for enlightening discussions. This work is supported in part by R01AI125982 (Y. Sun), R01DE024523 (Y. Sun), NYSDOH Rowley Breast Cancer Scientific Research Project DOH01-C33916GG-3450000 (Y. Sun), and R01CA241123 (S. Goodison).

Received July 25, 2019; revised August 30, 2019; accepted November 13, 2019; published first November 19, 2019.

References

1. Greaves M, Maley CC. Clonal evolution in cancer. *Nature* 2012;481:306–13.
2. Hanahan D, Weinberg RA. Hallmarks of cancer: the next generation. *Cell* 2011;144:646–74.
3. Davis A, Gao R, Navin N. Tumor evolution: linear, branching, neutral or punctuated? *Biochim Biophys Acta* 2017;1867:151–61.
4. Cancer Genome Atlas Network. Comprehensive molecular portraits of human breast tumours. *Nature* 2012;490:61–70.
5. Curtis C, Shah SP, Chin S-F, Turashvili G, Rueda OM, Dunning MJ, et al. The genomic and transcriptomic architecture of 2,000 breast tumours reveals novel subgroups. *Nature* 2012;486:346–52.
6. Ricketts CJ, De Cubas AA, Fan H, Smith CC, Lang M, Reznik E, et al. The cancer genome atlas comprehensive molecular characterization of renal cell carcinoma. *Cell Rep* 2018;23:313–26.
7. Sun Y, Yao J, Yang L, Chen R, Nowak N, Goodison S. Computational approach for deriving cancer progression roadmaps from static sample data. *Nucleic Acids Res* 2017;45:e69.
8. Cejalvo JM, de Duenas EM, Galvan P, Garcia-Recio S, Gasion OB, Pare L, et al. Intrinsic subtypes and gene expression profiles in primary and metastatic breast cancer. *Cancer Res* 2017;77:2213–21.
9. Johnson WE, Li C, Rabinovic A. Adjusting batch effects in microarray expression data using empirical Bayes methods. *Biostatistics* 2007;8:118–27.
10. Tibshirani R, Walther G, Hastie T. Estimating the number of clusters in a data set via the gap statistic. *J Royal Statist Soc* 2001;63:411–23.
11. Monti S, Tamayo P, Mesirov J, Golub T. Consensus clustering: a resampling-based method for class discovery and visualization of gene expression microarray data. *Mach Learn* 2003;52:91–118.
12. Sørlie T, Perou CM, Tibshirani R, Aas T, Geisler S, Johnsen H, et al. Gene expression patterns of breast carcinomas distinguish tumor subclasses with clinical implications. *Proc Natl Acad Sci U S A* 2001;98:10869–74.
13. Parker JS, Mullins M, Cheang MC, Leung S, Voduc D, Vickery T, et al. Supervised risk predictor of breast cancer based on intrinsic subtypes. *J Clin Oncol* 2009;27:1160–7.
14. Goldhirsch A, Wood WC, Coates AS, Gelber RD, Thurlimann B, Senn H-J, et al. Strategies for subtypes – dealing with the diversity of breast cancer: highlights of the St Gallen International Expert Consensus on the primary therapy of early breast cancer 2011. *Ann Oncol* 2011;22:1736–47.
15. Creighton CJ. The molecular profile of luminal B breast cancer. *Biologics* 2012;6:289–97.
16. Mackay A, Weigelt B, Grigoriadis A, Kreike B, Natrajan R, A'Hern R, et al. Microarray-based class discovery for molecular classification of breast cancer: analysis of interobserver agreement. *J Nat Cancer Inst* 2011;103:662–73.
17. Nowell PC. The clonal evolution of tumor cell populations. *Science* 1976;194:23–8.
18. Beerenwinkel N, Schwarz RF, Gerstung M, Markowitz F. Cancer evolution: mathematical models and computational inference. *Syst Biol* 2015;64:e1–e25.
19. Anderson WF, Rosenberg PS, Prat A, Perou CM, Sherman ME. How many etiological subtypes of breast cancer: two, three, four, or more? *J Natl Cancer Inst* 2014;106:165.
20. Berger AC, Korkut A, Kanchi RS, Hegde AM, Lenoir W, Liu W, et al. A comprehensive pan-cancer molecular study of gynecologic and breast cancers. *Cancer Cell* 2018;33:690–705.