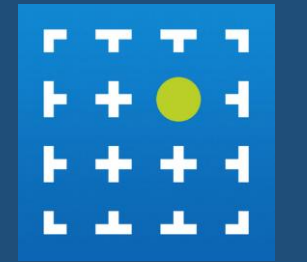


# Machine learning approach to personalized medicine in breast cancer patients: development of data-driven, personalized, causal modeling through identification and understanding of optimal treatments for predicting better disease outcomes



Kaplan HG<sup>1</sup>, Berry AB<sup>1</sup>, Rinn KJ<sup>1</sup>, Ellis ED<sup>1</sup>, Birchfield GR<sup>1</sup>, Wahl TA<sup>1</sup>, Liu X<sup>1</sup>, Tameishi M<sup>1</sup>, Beatty JD<sup>1</sup>, Dawson PL<sup>1</sup>, Mehta VK<sup>1</sup>, Holman A<sup>1</sup>, Atwood MK<sup>1</sup>, Alexander S<sup>1</sup>, Bonham C<sup>1</sup>, Summers L<sup>1</sup>, Khalil I<sup>2</sup>, Hayete B<sup>2</sup>, Wuest D<sup>2</sup>, Zheng W<sup>2</sup>, Liu Y<sup>2</sup>, Wang X<sup>2</sup>, Brown TD<sup>1</sup>.  
<sup>1</sup>Swedish Cancer Institute, Swedish Medical Center, Seattle. <sup>2</sup>GNS Healthcare, Cambridge.

## BACKGROUND

In the era of personalized medicine, a major challenge is harnessing longitudinal data across the cancer care continuum, which includes multi-modal data sets of biological, molecular, and clinical information about patients (pts) and their tumors. There is a growing need for new computing analytics, such as machine learning - an important tool in healthcare bio-informatics. We report our approach to building cancer disease models in an unbiased manner through utilization of a causal machine learning and simulation platform.

## METHODS

The Swedish Cancer Institute (SCI) Personalized Medicine Research Program (PMRP) is a prospective registration protocol with the objective of establishing a centralized longitudinal, molecular, and phenotypic data repository. Since 2014, over 1,070 pts have been enrolled, having undergone next generation sequencing (NGS) profiling of their tumors. Of these pts, we identified 100 breast cancer pts who also have detailed longitudinal clinical annotation within our SCI Breast Cancer Registry (BR). The BR was established in 1990 and has over 17,000 cases. BR follows all newly diagnosed breast cancer pts treated at SCI. Detailed information is collected on demographics, clinical presentation, tumor characteristics, staging, treatment and follow-up. The biomarker distribution and time in years from breast cancer diagnosis and NGS test for 100 pts are summarized in Fig. 1 and Fig. 2.

All de-identified data, variables, and data points in the multi-modal data types are integrated into normalized data frames. A reverse engineering approach, via the Reverse Engineering and Forward Simulation (REFS) platform, is being utilized, focusing on discovering the complex causal mechanisms that determine which therapies will produce the best outcomes for an individual pt. This method goes beyond traditional approaches that rely on data correlations to match treatments to pts. The breast cancer causal model uncovers many of the possible combinations of causal relationships that drive outcomes and enables "what if?" simulations of a variety of interventions, across pts, to determine optimal therapies. Performance metrics and model robustness will be explored using a stratified, n-fold (e.g. 10-fold) cross-validation procedure, which is designed to provide an unbiased estimate of model generalization to new observations. (Fig. 3)

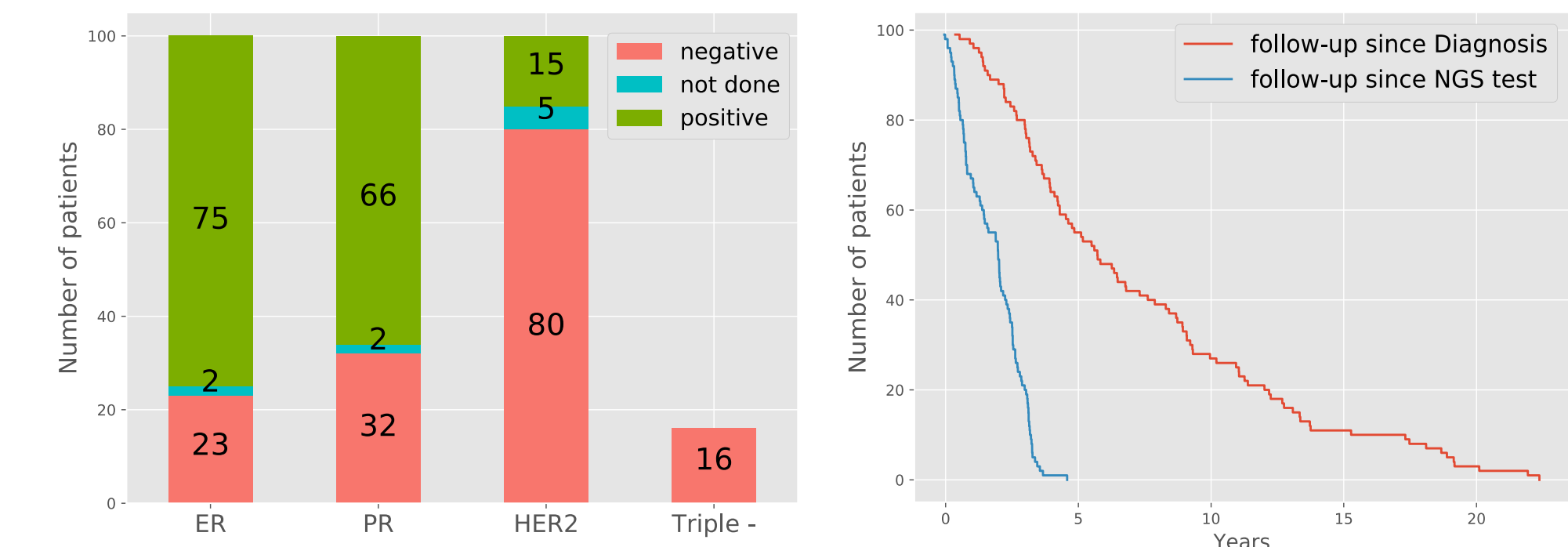


Fig 1. Pts Characteristics - Biomarkers

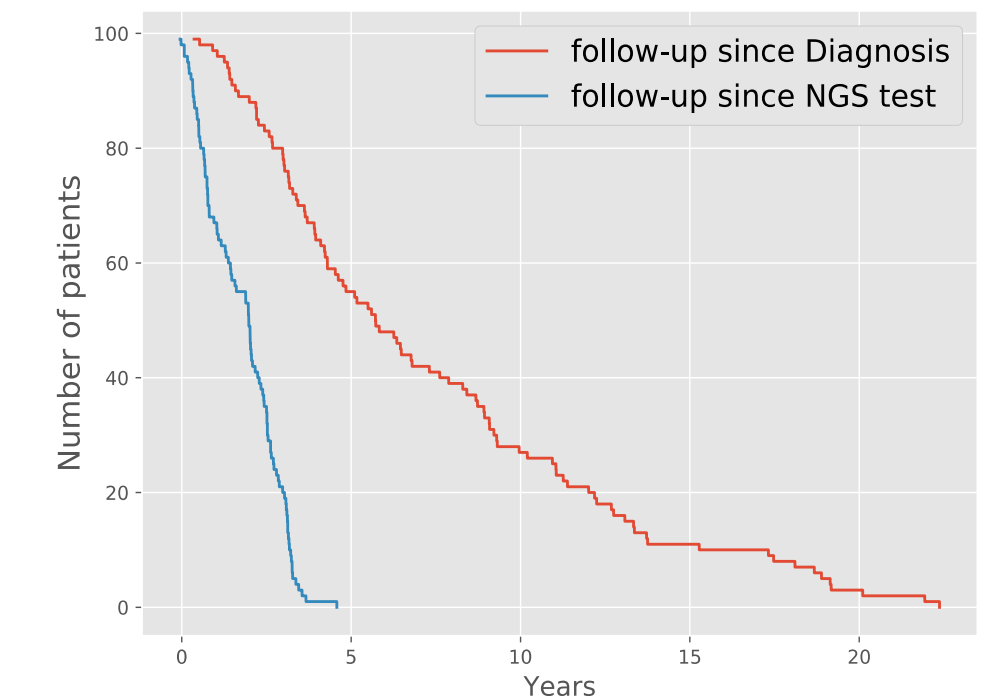
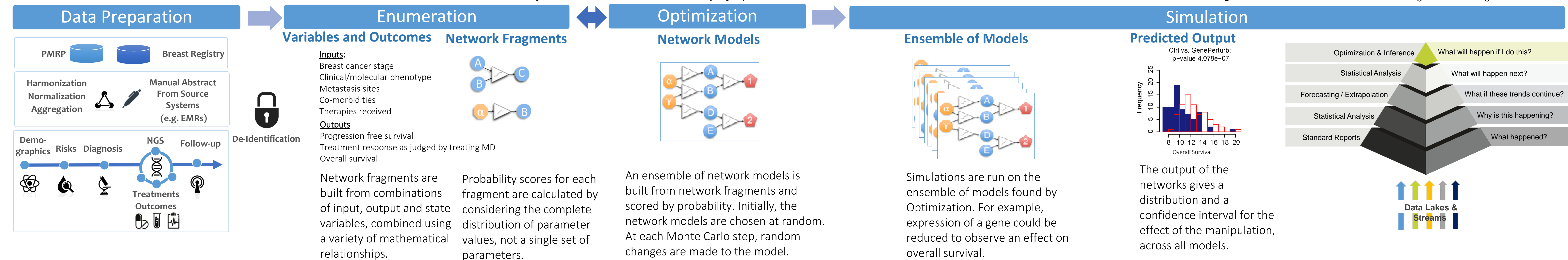


Fig 2. Time from diagnosis and NGS tests

Fig. 3 Breast Cancer Model: Identifying Optimal Treatments in Patients



## RESULTS

The causal model and simulations can elevate the providers' abilities to: better understand treatment responses based on pts' unique clinical data and mutational statuses; study different treatment options to optimize management; and understand the complex interactions among variables that lead to a range of treatment outcomes.

Challenges experienced as of today:

- Data access and availability (e.g. system changes over time; lack of access to outside records due to care transfer; lack of EMR automatic download)
- Lack of complete discrete clinical data and reliance on manual curation
- Aggregation of different types of data (e.g. time-series data; static data)
- Prioritizing/balancing granularity vs. breadth of data to guide manual abstraction efforts
- Comparability of data with classification changes (e.g. cancer staging and Her2/FISH)
- Retrospective analysis of cancer biomarkers with phenotype/genotype changes in longitudinal data
- Harmonization and normalization of multiple databases (e.g. coding and data classification)
- Timing of NGS testing vs. clinical application

## CONCLUSIONS

Machine learning could potentially provide novel insights into personalized medicine. There are challenges to creating a large enough clean data set that when analyzed will produce results that can be confidently used to inform patient care.

Knowledge generated from the simulations of the disease model can potentially streamline and support the clinical decision-making process, to include molecular tumor board deliberations, and ultimately assist providers in arriving at optimal treatment recommendations for pts.

## ACKNOWLEDGEMENTS

We thank CellNetix Pathology and Laboratories for providing majority of tumor genomic profiling.