# PREDICTING THE EFFECTS OF WAVE GENES ON CANCER SEVERITY AND OUTCOMES

Ekaterina Cole, SUNY Oswego, Biomedical Informatics MS

Joshua Harkness, SUNY Oswego, Information Science BA

Seham Al-Masri, SUNY ESF, Biotechnology BS

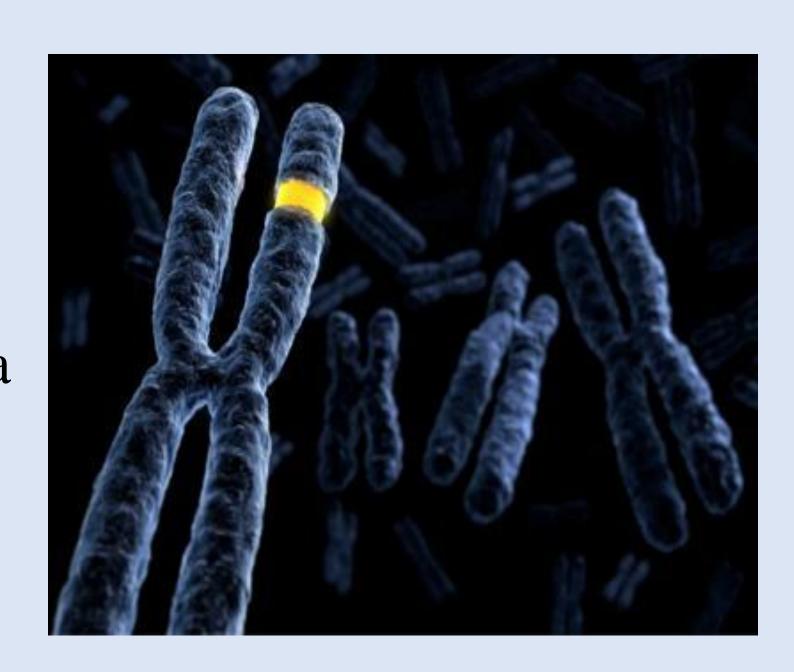Dr. Isabelle Bichindaritz, SUNY Oswego, Director of Biomedical Informatics

**Abstract:**
The WAVE regulatory complex has been identified as one of the most important actors contributing to breast tumor progression. The WAVE is represented by a set of 11 genes, namely, ABI1 ABI2 ABI3 BRK1 CYFIP1 CYFIP2 NCKAP1 NCKAP1L WASF1 WASF2 WASF3, functioning in a synergistic manner. The main objective was to define the characteristics of the WAVE and determine its role in breast cancer as a phenomenon. A multi-systemic dataset, available in the public domain, was obtained, structured and analyzed. This project implied a collaboration with Dr. Kotula, an oncologist at Upstate Medical Center.

**Introduction and Stats:**
Breast cancer is one of the most omnipresent diseases in today's US Healthcare system. According to the American Cancer Society statistics (2017), in **2017** alone, an estimated **40,610** women will **die** from breast cancer while approximately **316,120** women will be **diagnosed**; a fourth of which will be non-invasive. In addition, breast cancer is reported to be the second leading cause of death for women behind lung cancer. The chance a woman will die from breast cancer is **1 in 37**. These data emphasize the importance of most profound understanding of factors that trigger breast cancer and contribute to its development.

**Project Milestones and Methods:**
The project encompassed examining the professional domain of breast cancer and identifying the existing research. Meta-analysis of literature was expected to provide basis for data understanding.
The breast cancer data, available in public domain (TCGA), was obtained and merged with the help of a specifically designed R script. The dataset, including both clinical and genomic data, was analyzed and the full data dictionary was completed. The WAVE complex was extracted from the dataset and combined with a set of variables, descriptive of

- **Cancer stage;**
- **Occurrence of relapse;**
- **Presence of metastasis.**

**The R Script:**
The R script written merges clinical data with genomic data. The datasets were merged based on a shared barcode. Editions were made to the clinical dataset for the barcodes to match exactly, due to minor differences. This ensured that the proper barcodes were given additional variables from clinical, without any repeated values.

An outer join was used to ensure no variables were lost while still merging identical barcodes together. The merge introduced an additional 3719 variables.
The WAVE characteristics were summarized using a range of statistical tools, including Excel, R, and SPSS. The associations analysis, along with dataset enrichment through cBioPortal was performed.

**Project Achievements:**
- A merged dataset of clinical + genes data, with a designed R script, and full data dictionary for the dataset were created;
- WAVE complex extracted and variables, definitive of main set characteristics were identified;
- The characteristics of the WAVE were summarized and the associations have been examine:
    - the highest values of gene expressions are associated with the higher cancer stages (X, IIIB), and the lower values, are associated with the stages II and III.
    - genes of BRK, CYFIP and NCKAP families appear to be among the most frequently expressed genes for cases with both relapse occurrence and higher stages;
    - significant differences in WAVE gene expression levels have been observed in connection with aggregated patients' replapse status;
    - cBioPortal Enrichment analysis has shown that Hormonal genes of BRCA and HER families have been found to be correlated with the WAVE sequence in terms of genetic alterations.

| Name of the Dataset Feature | Characteristics |
|---|---|
| Size of the Overall Dataset Analyzed (Clinical+Genes) | 111025 rows x 2242 columns (1.5GB) |
| Complete Number of Genes | 21661 gene (5 metrics per each) |
| WAVE genes | 11 genes extracted |
| Number of patients | 2238 patients with unique TCGA barcodes |
| Main Patients' Characteristics | Patient Barcode, Vital Status, New tumor event; Days to death, Tumor status; Cancer Stage; Drug/ therapy info; follow-ups. |
| Characteristics, Vital to Data Analysis | Patient's vital status; cancer stage; relapse status; metastasis status; patient's new tumor events and follow-up condition statuses. |
| Origin of the Data | NIH The Cancer Genome Atlas (TCGA) |