

Published: November 30, 2023

Citation: Linder, K., et al., 2023. Prediction of Ovarian Cancer with Deep Machine Learning and Alternative Splicing. Medical Research Archives, [online] 11(11). <https://doi.org/10.18103/mra.v11i11.4602>

Copyright: © 2023 European Society of Medicine. This is an open-access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

DOI: <https://doi.org/10.18103/mra.v11i11.4602>

ISSN: 2375-1924

RESEARCH ARTICLE

Prediction of Ovarian Cancer with Deep Machine Learning and Alternative Splicing

Katharine Linder¹, Rachel Watson¹, Keely Ulmer¹, David Bender¹, Michael J Goodheart¹, Eric Devor¹, Jesus Gonzalez Bosquet^{1*}

¹University of Iowa Hospitals and Clinics Department of Obstetrics and Gynecology, Iowa City, Iowa.

*jesus-gonzalezbosquet@uiowa.edu

ABSTRACT

Objective: Early detection of ovarian cancer could lead to improved survival rates, however no method has reliably been able to predict ovarian cancer. The aim of this study is to determine if processing alternative splicing data from high grade serous ovarian cancer patients using machine learning analytics will discriminate high grade serous ovarian cancer from normal fallopian tube samples. The ultimate goal would be to have a model that can predict high grade serous ovarian cancer with a blood test.

Methods: This is a case-control study of patients with confirmed high grade serous ovarian cancer and those undergoing salpingectomy for benign indications. RNA-sequencing was performed on all samples. RNA-sequence data was then put into Deep-learning augmented RNA-seq analysis of transcript splicing software suite. Deep-learning augmented RNA-seq analysis of transcript splicing created a model of differential alternative splicing aimed to discriminate between high grade serous ovarian cancer and normal fallopian tube. DEXSeq analysis was used to determine exon-based expression. Initial results with both analytics were then modelled with multivariate lasso regression to create prediction models (performance determined by area under the curve and 95% CI). Models created were the validated using The Cancer Genome Atlas data sets.

Results: One hundred and twelve high grade serous ovarian cancer and 12 benign samples were successfully sequenced. Deep-learning augmented RNA-sequencing analysis of transcript splicing identified 998 unique differentially expressed exons between high grade serous ovarian cancer and controls. Multivariate lasso regression analysis identified several exons that predicted high grade serous ovarian cancer with high performance. Specifically, ENSG00000182512:E001 from gene GLRX5 was highly predictive of high grade serous ovarian cancer with an area under the curve of 100%.

Conclusions: Application of machine learning analytics to exon differential expression, most likely due to alternative splicing, predicted high grade serous ovarian cancer with high performance. These results were validated in an independent dataset of cases and controls. Differential exon expression from cell-free RNA potentially could be used for early diagnosis of high grade serous ovarian cancer.

Keywords: ovarian cancer, alternative splicing, machine learning.

Introduction

There are estimated to be less than 20,000 new cases of ovarian cancer in the United States in 2022, but it is the fifth leading cause of cancer deaths in women¹. Of the subtypes of ovarian cancer, high-grade serous carcinoma (HGSC) is the most common and encompasses cancers originating from the fallopian tube, ovary, and peritoneum. Ovarian cancer is most frequently diagnosed at advanced stages, and the 5-year survival rate for patients with ovarian cancer is less than 50%. However, when diagnosed at an early stage it is greater than 90%². Early detection and treatment of ovarian cancer has been shown to decrease mortality, however, there is a lack of screening methods to achieve early detection³.

Patients with ovarian cancer will typically present with a pelvic mass, however, the differential diagnosis for pelvic mass is broad and preoperative diagnosis of pelvic masses can be difficult. While cancer antigen 125 (CA-125) and transvaginal ultrasound (TVUS) are critical for the initial assessment of adnexal mass, each lacks sensitivity and specificity to be utilized as standard screening tests. CA-125 and TVUS have been utilized in large studies but failed to show a significant mortality benefit and illustrated the potential for harm for patients with benign masses undergoing unnecessary surgical intervention^{4,5}. Studies have shown that patients with gynecologic cancers have better outcomes when treated by gynecologic oncologists⁶⁻⁸. Improved screening methods would lead to better healthcare utilization and patient outcomes by assessing which patients with adnexal masses should be evaluated by gynecologic oncologists.

There is emerging evidence for utilizing cell free DNA (cfDNA) and circulating tumor DNA (ctDNA) in a patient's plasma⁹. Circulating tumor DNA refers to the tumor DNA that is distinct from the patient DNA present in the plasma. This has become known as a "liquid biopsy." Studies have shown the potential for cfDNA and ctDNA to be utilized in the diagnosis and treatment strategies for several cancers including pancreatic, colorectal, lung, and breast¹⁰⁻¹³. We sought to create a machine learning (ML) prediction model analyzing RNA sequencing that could be used as a "liquid biopsy" and could diagnose HGSC patients at earlier stages.

Deep ML is a means of utilizing artificial intelligence to enhance our ability to process large amounts of data and make predictive models. Deep ML models use algorithms to analyze patterns in data sets, then using these patterns the model trains itself to make predictions on new data sets¹⁴. It has been used in oncology with varied applications, such as to classify disease genetic variations in genomes and identify RNA sequences^{15,16}. A prior study developed Deep-Learning Augmented RNA-seq analysis of Transcript Splicing, or DARTS. This program involves a "deep neural network model that predicts differential alternative splicing between two conditions based on exon-specific sequence features and sample specific regulatory features; and a Bayesian Hypothesis (BHT) statistical model that infers differential alternative splicing by integrating empirical evidence in a specific RNA-seq dataset with prior probability of differential alternative splicing¹⁷." We will use DARTS, a ML based method, to create a prediction model that could detect HGSC in blood.

DEXseq is another method that has been used to test for differences in the usage of exons in RNA sequencing samples and, consequently, alternative splicing in biological samples. DEXseq has the ability to control for false positives by taking into account biological variability¹⁸. We will also use this method to assess differential exon expression to compare and complement DARTS method.

We utilized a large biobank of ovarian cancer specimens and normal fallopian tubes at our institution to analyze RNA splicing. We hypothesize that alternative splicing analyzed with ML will discriminate HGSC from normal fallopian tube samples. Our objective is to create a model to identify HGSC using two different methods for differential exon expression and alternative splicing, deep ML framework (DARTS) and DEXSeq methods. These models were further validated with independent datasets, TCGA, and in ML analytical platforms.

Methods

We performed a single-institution retrospective case-control study using HGSC tumor specimens obtained at the time of cytoreductive surgery and benign fallopian tube specimens. RNA was isolated from all specimens and RNA sequencing (RNA-seq) was performed.

SPECIMEN ACQUISITION

High grade serous ovarian cancer tissue samples were obtained from the Department of Obstetrics and Gynecology Gynecologic Oncology Biobank (IRB, ID no. 200209010), part of the Women's Health Tissue Repository (WHTR, IRB, ID no. 201804817) at the University of Iowa (UI). All tissues archived in

the Gynecologic Oncology Biobank (herein termed Biobank) were originally obtained from adult patients under informed consent in accordance with University of Iowa Institutional Review Board (IRB) guidelines. Tumor samples were collected, reviewed by a board-certified pathologist, flash-frozen, and then the diagnosis was confirmed in paraffin at the time of initial surgery. All experimental protocols were approved by the University of Iowa Biomedical IRB-01.

Fallopian tube samples from patients undergoing gynecologic procedures were obtained from patients with no family history of cancer beside squamous cell carcinoma of the skin and who were undergoing salpingectomy for benign indications. (IRB, ID no. 201202714). These samples were similarly obtained from adult patients under informed consent in accordance with the University of Iowa IRB guidelines. RNA from both the fallopian tubes and HGSC specimens had already been extracted and purified in a previous study.

Two-hundred and fifty-three patients with advanced (stage II and IV) or recurrent HGSC were identified from the database. Of these 193 had available flash frozen tissue for RNA isolation. One-hundred and twelve of the samples had good quality RNA successfully processed for RNA-seq. Of the 20 benign fallopian tube samples obtained, 12 patients had good quality RNA that was successfully processed for RN-seq (**Figure 1A**). All samples were obtained and processed for previous studies¹⁹.

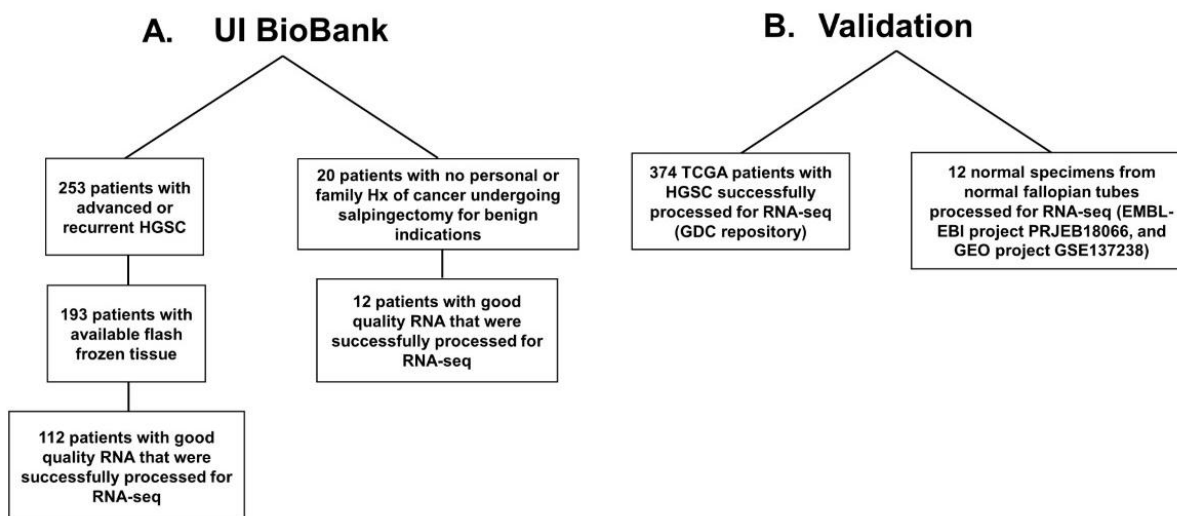


Figure 1: Flow chart of specimens used for analysis. A. Specimens processed from UI Biobank: 112 cases and 12 controls (normal fallopian tubes). B. Publicly available samples for validation: TCGA HSC specimens, and normal fallopian tubes from EMBL-EBI repository (project number PRJEB18066) and GEO accession number GSE137238.

RIBONUCLEIC ACID SEQUENCING

Ribonucleic acid was isolated from HGSC and benign fallopian tube samples. Methods for RNA sequencing and processing have been described elsewhere¹⁹. Total cellular RNA was extracted from tissue with the mirVana (Thermo Fisher, Waltham, USA) RNA purification kit. The RNA yield and quality were assessed with Trinean Dropsense 16 spectrophotometer and Agilent Model 2100 bioanalyzer. RNA quality was determined to be adequate if the sample had an RNA integrity number (RIN) of 7.0 or greater. Samples that were of adequate quality were then sequenced. 500ng of RNA was quantified by Qubit measurement (Thermo Fisher). RNA was then converted to cDNA and ligated to sequencing adaptors with Illumina TriSeq stranded total RNA library preparation (Illumina, San Diego, CA, USA). Sequencing was then carried out on the Illumina HiSeq 4000 genome sequencing platform using 150 bp paired-end SBS chemistry. All sequencing was performed in the Genome Facility of the University of Iowa Institute of Human Genetics (IIHG).

DEEP-LEARNING AUGMENTED RNA-SEQ ANALYSIS OF TRANSCRIPT SPLICING ANALYSIS

Deep-learning augmented RNA-seq analysis of transcript splicing (DARTS) is a method for predicting alternative splicing patterns¹⁷. It uses a bayesian hypothesis testing (BHT) model paired with a deep neural network (DNN) to predict differential alternative splicing. The DARTS BHT initially analyzes large-scale RNA-seq data to generate training labels for differential or unchanged splicing events. This is used to train the DNN. Then the trained DARTS DNN was used to predict alternative splicing and differential exon expression in our data set. The prediction, along with observed RNA-seq read counts, are incorporated as an informative prior by the DARTS BHT to perform deep learning augmented splicing analysis.

DEXSEQ ANALYSIS

DEXSeq analysis was used to determine exon-based expression¹⁸. BAM files from RNAseq were used to determine expression of single

exons. The DESeq2 R package was used to normalize and transform data to determine the differences between HGSC samples and benign fallopian tube samples²⁰ (Supplementary Figure S1 and S2). The differential expression analysis in DESeq2 package uses a generalized linear model to assess difference in log₂-transformed counts. A p-value (Wald test p-value) adjusted with false discovery rate (FDR) was considered significant when < 0.0001 to account for multiple comparisons²¹.

CREATION OF PREDICTION MODELS OF HIGH GRADE SEROUS OVARIAN CANCER WITH MULTIVARIATE MODELS

Significant variables from the univariate analysis - ANOVA for the DEXSeq analysis and DARTS analysis - were then incorporated into two different multivariate lasso regression prediction models of ovarian cancer (HGSC), one for each method. Multivariate prediction models were fit with lasso (least absolute shrinkage and selection operator) as implemented in the glmnet R package²², and detailed previously²¹. Performances of prediction models were measured by the area under the receiver operating characteristics curve (AUC) and their 95% confidence interval (CI), estimated with 1,000 replicates of ten-fold cross-validation to avoid over-fitting. Bias-corrected and accelerated bootstrap CI's were computed for each model. AUC of 0.5 indicates no predictive ability of a model and 1.0 represents perfect predictive performance.

VALIDATION OF PREDICTIVE MODELS WITH THE CANCER GENOME ATLAS DATABASE AND MACHINE LEARNING ANALYTICS

High grade serous ovarian cancer data from The Cancer Genome Atlas (TCGA) was used

to validate created predictive models²³. We included publicly available controls from EMBL-EBI project PRJEB18066 and GEO project GSE137238 databases (Figure 1B). After permission was granted to access controlled data by the Genomic Data Commons (GDC) Data Portal (dbGaP# 29868), TCGA HGSC BAM files from RNA-seq experiments aligned to the human reference genome (version hg38) with the STAR suite were downloaded in their original format. DEXSeq analysis to determine exon-based expression was performed as described previously for both, TCGA cases and GEO and EMBL-EBI controls. Then selected features used in the multivariate prediction lasso model were extracted also from TCGA single exon expression dataset and used for prediction model validation of both, DARTS and DEXSeq analyses.

Machine learning prediction models with TensorFlow platform were performed with those features more informative for prediction of HGSC in previous analysis. Initial validation of prediction models with ML included all selected exons from initial univariate and multivariate analyses for both DARTS and DEXSeq. Then we used ML platforms to validate models performed with those exons selected in TCGA validation analysis. The final goal was to identify and validate the simplest, most accurate and robust model that could predict HGSC.

Results

We extracted count expression information from 645,243 exons from over 60,000 different transcripts in our database, including genes, long noncoding RNAs (lncRNAs), MIRs and other transcription units. Differential

exon expression was performed with both DARTS and DEXSeq analysis (Figure 2). Deep-learning augmented RNA-seq analysis of transcript splicing identified 998 unique differentially expressed exons between HGSC

and controls (Figure 2A, upper panel). A prediction model with all these exons had an AUC of 91%. Differential exon expression with DEXSeq, identified 1,399 significant exons ($p < 0,001$).

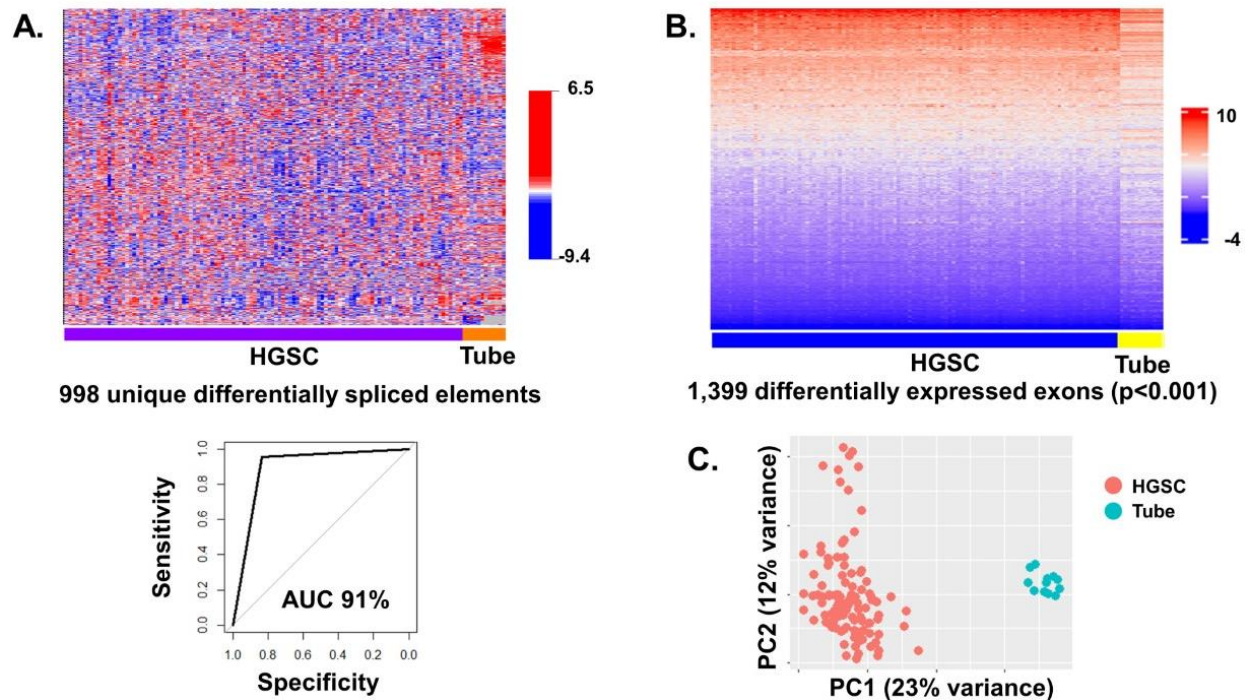


Figure 2: Analysis of differential exon expression and alternative splicing by different methods: **A.** DARTS: first a deep neural network (DNN) model predicts differential alternative splicing based on exon-specific sequence features and sample-specific features (upper panel); then, a Bayesian hypothesis testing (BHT) statistical model infers differential alternative splicing by integrating empirical evidence in a specific RNA-seq dataset with prior probability of differential alternative splicing. Prediction model represented by the ROC curve (lower panel). **B.** Differential exon expression by DEXSeq analysis. Heatmap with resultant exons. **C.** Principal component (PC) analysis of differential exon expression with DEXSeq. The figure represents a plot of the two main PCs.

CREATION OF MULTIVARIATE PREDICTION MODELS OF HIGH GRADE SEROUS OVARIAN CANCER

A multivariate lasso regression prediction model using all features selected in the univariate analysis after DEXSeq assessment (N=1,399) identified a unique exon, ENSG00000182512:E001 (within gene *GLRX5*) that predicted HGSC with an AUC of 100% (Figure 3A, upper panel). The multivariate lasso regression using the results from the DARTS analysis identified three

exons: ENSG00000050130:E013 (within gene *JKAMP*), ENSG00000135597:E014 (in gene *REPS1*), and ENSG00000175061:E013 (in lncRNA *FAM211A-AS1*, or *SNHG29*), which predicted HGSC with an AUC of 100% (Figure 3A, lower panel). When assessing the relative expression between of these 4 exons between HGSC and tube samples, only ENSG00000182512:E001 increased risk for HGSC, the other three were decreased, and protective of HGSC (Figure 3B).

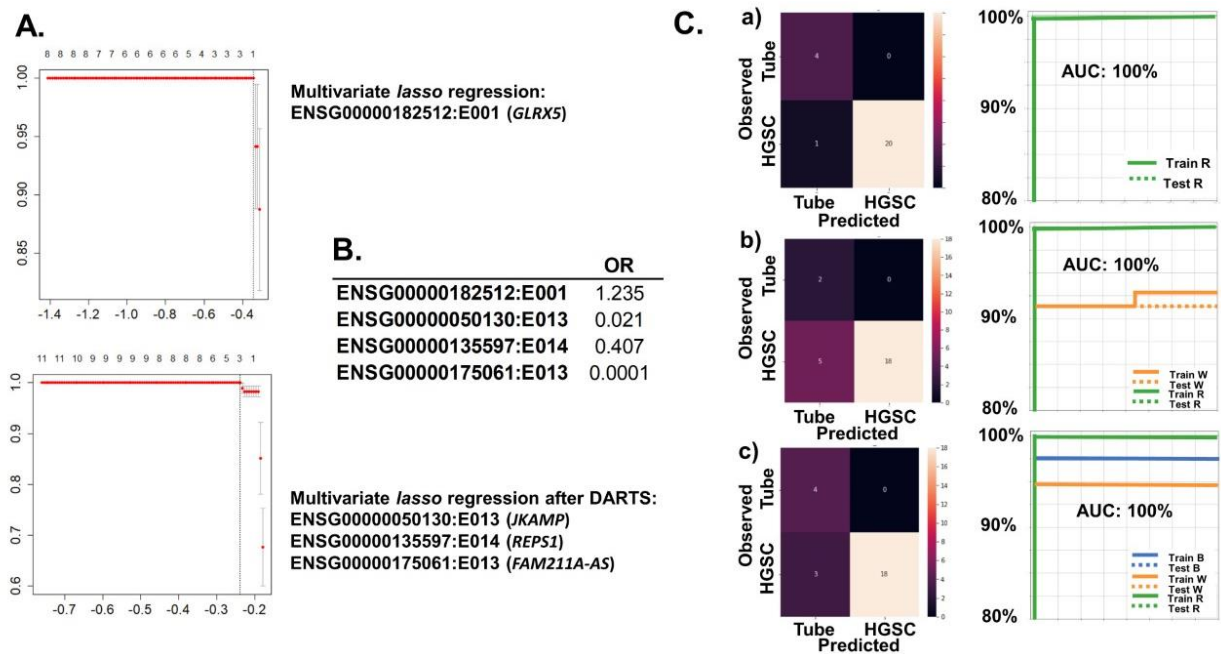


Figure 3: Multivariate lasso regression analysis and machine learning validation with TensorFlow.

A. Multivariate lasso regression analysis to predict HGSC with differential exon expression: In the upper panel, multivariate lasso regression prediction model using results from univariate regression analysis after DEXSeq ($p < 0.001$): one differentially expressed exon, ENSG00000182512:E001 (*GLRX5*) predicted HGSC with an AUC of 100%; the inferior panel shows the multivariate lasso regression prediction model using results from the machine learning augmented method DARTS: three exons, ENSG00000050130:E013 (*JKAMP*), ENSG00000135597:E014 (*REPS1*), and ENSG00000175061:E013 (*FAM211A-AS*), predicted HGSC with an AUC of 100%. **B.** Exons odds ratio expression with respect to normal Fallopian tube: only ENSG00000182512:E001 (*GLRX5*) expression was elevated in HGSC vs tube; the expression of the other 3 exons were lower in HGSC samples. **C.** Validation of HGSC prediction models with machine learning using TensorFlow: **a)** Model using all differentially expressed exon in DEXSeq univariate analysis (N=1,399): the left panel shows the confusion matrix representing the observed versus the predicted values; the right panel represents the ROC graphic including a model accounting for unbalanced samples: Train R: results of unbalanced (or re-sampling) model training; Test R: results of re-sampling model testing. **b)** Model using all 998 unique differentially spliced elements in the DARTS analysis (N=998): the left panel shows the confusion matrix representing the observed versus the predicted values; the right panel represents the ROC graphic including: 1) models accounting for weights of the outcome: Train W: results of weighted model training; Test W: results of weighted model testing; 2) models accounting for unbalanced samples: Train R: results of unbalanced (or re-sampling) model training; Test R: results of re-sampling model testing. **c)** Model using all 4 exons found to be significant in both multivariate lasso analyses: the left panel shows the confusion matrix representing the observed versus the predicted values; the right panel represents the ROC graphic including: 1) basic model: Train B: results of basic model training; Test B: results of basic model testing; 2) models accounting for weights of the outcome: Train W: results of weighted model training; Test W: results of weighted model testing; 3) models accounting for unbalanced samples: Train R: results of unbalanced (or re-sampling) model training; Test R: results of re-sampling model testing.

VALIDATION OF PREDICTION MODELS

Initially, these multivariate prediction models were validated in a different ML analytical platform using *TensorFlow*. First, we validated the models including features selected after the univariate analysis with DEXSeq and the

initial DARTS assessment. The performance of prediction models with 1,399 exon from the DEXSeq analysis (Figure 3C.a) and 998 exons selected by DARTS (Figure 3C.b) had excellent performances, with AUCs of 100%. A third model with only the 4 resulting exons

after both multivariate lasso regressions (for DEXSeq and DARTS results), resulted in another model with excellent performance also, AUC of 100% (Figure 3C.c).

Then, these models were validated in an independent dataset. After downloading, homogenizing, determining exon expression, we normalized and log2 transformed TCGA and control data from EMBL-EBI and GEO databases. Out of the initial 1,399 exons selected after univariate analyses with DEXSeq, 1,171 were also detected in

TCGA+controls databases. A new multivariate lasso regression was done with those 1,171 exons in UI dataset, and then the model was applied to TCGA+controls dataset using *pROC* (a R package) in data (Figure 4A). The same unique exon, ENSG00000182512:E001 (in gene *GLRX5*) predicted HGSC in UI data with an AUC of 100%, and when the model was applied to TCGA data, the performance of the validated model was excellent, with an AUC of 97% (Figure 4C, upper panel).

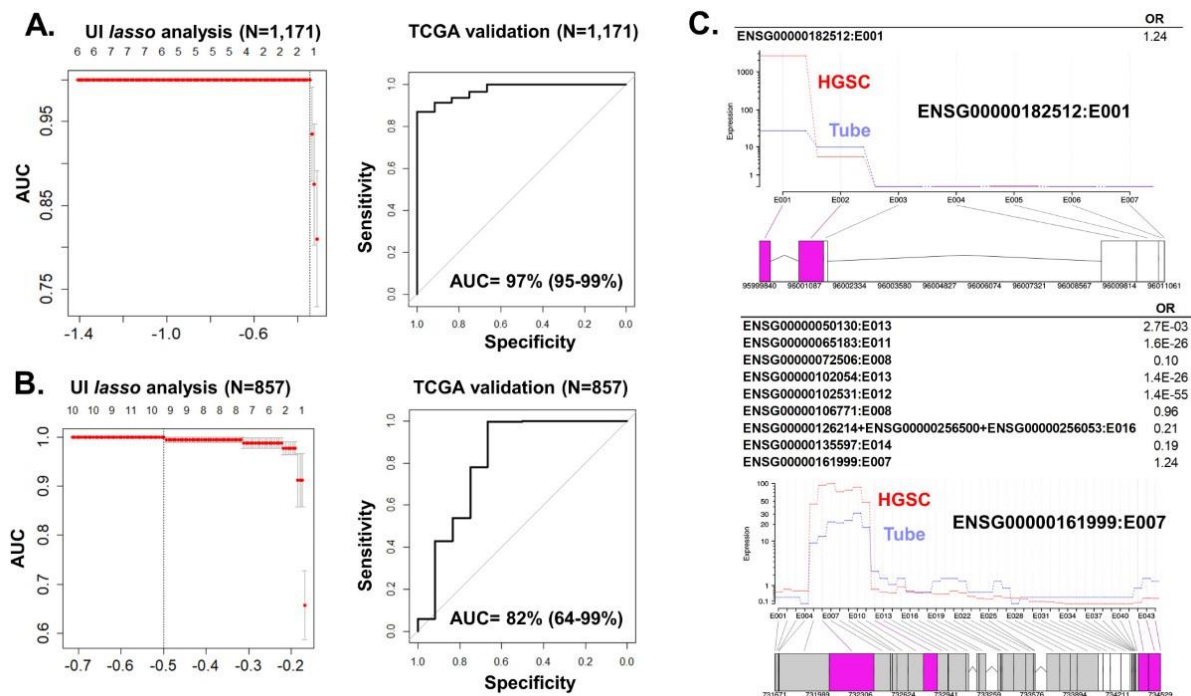


Figure 4: Validation of differential exon expression with TCGA data.

A. Validation of HGSC prediction with differential exon expression of DEXSeq analysis: the left panel, shows the prediction model using exon expression significant in the UI cohort, that is also present in TCGA cohort (N=1,171): again, one differentially expressed exon, ENSG00000182512:E001 (*GLRX5*) predicted HGSC with and AUC of 100%. The right panel shows the validation of the prediction model in TCGA data, with an excellent performance of 97% (CI, 95-99%).

B. Validation of HGSC prediction with differential exon expression regression prediction model using results from the machine learning augmented method DARTS: the left panel, shows the prediction model using exon expression also present in TCGA cohort (N=857): 9 differentially expressed exons predicted HGSC with and AUC of 100%. The right panel shows the validation of the prediction model in TCGA data, with a very good performance of 82% (CI, 64-99%).

C. Detailed level of exon expression in 2 of the exons. Both, ENSG00000182512:E001 and ENSG00000161999:E007, had higher expression in HGSC samples (red) than in tubal samples (blue). Comparison of other exons in the same genes are also represented.

The strengths of this study are the use of data from a clinically well annotated database of patients with HGSC with comprehensive genomic information and an adequate number of controls. Additionally, prediction models were validated in an independent comprehensive dataset (TCGA). Further, data were analyzed with two separate analytical platforms. The limitations of this study are that it is retrospective, and therefore is at risk of bias. The UI database is limited by minimal diversity in the sample given the demographics of Iowa, with the general population made up of 90.1% white individuals⁴⁰. Future studies using these models prospectively and in larger data sets are needed to test their clinical utility.

Conclusion

Our study created a validated model for identifying HGSC by utilizing deep machine learning framework (DARTS) and DEXseq to identify differences in exon expression and alternative splicing in tumor cells when compared to benign fallopian tube cells. In the future, this model could be utilized to differentiate tumor cells in blood samples as a means for early and non-invasive diagnosis of HGSC, which could ultimately lead to a reduction in mortality.

Conflict of Interest Statement:

The authors have no conflicts of interest to disclose.

Funding Statement:

None

Acknowledgement Statement:

Funding: This work was supported in part by the NIH grant R01 CA99908 and R01 CA184101 to Kimberly K. Leslie, and the basic research fund from the Department of Obstetrics & Gynecology at the University of Iowa. Also, was supported in part by the American Association of Obstetricians and Gynecologists Foundation (AAOGF) Bridge Funding Award.

