

# Machine & Deep Learning from Longitudinal EHRs and Genetic Data for Low CVD Risk Prediction

Matthew Hancock<sup>1,2</sup>  
Mentor: Wei-Qi Wei, MD, PhD<sup>2</sup>

<sup>1</sup>Vanderbilt University, Department of Mathematics  
<sup>2</sup>Vanderbilt University Medical Center, Department of Biomedical Informatics

**OBJECTIVE:**  
*Our study sought to use machine & deep learning to leverage recent advances in high-throughput genotyping & phenotyping to learn from longitudinal EHR, genomic, and expression data to predict CVD outcome in the conventionally low-risk population.*

**Introduction**  
Cardiovascular disease (CVD) is the leading cause of morbidity/mortality globally, responsible for 1/3 of deaths.

The **Framingham Risk Score (FRS)** is the gold standard for 10-year CVD prediction, but has key limitations:

- Predicting a complex disease using only a **small number of risk factors**
- Attempting to analyze disease outcome over time using **cross-sectional data**
- Placing the **highest correlation coefficient on age**

15-20% of first CVD-event patients had one or zero conventional CVD risk factors.

**Study Design**  
Individuals for the study were selected from the VUMC Synthetic Derivative (SD), an anonymized copy of the whole VUMC EHR (>3 million unique individuals).

Study cohort	
Age: 18 – 89	
Case: ≥ 1 CVD diagnosis code(s) ; Control: no CVD code	
1:8 case-control match (age, sex, race, and T <sub>0</sub> [first diagnosis of CVD])	
EHR data within 5-year window prior to T <sub>0</sub> (6-month prior to diagnosis to avoid confounding issue)	
FRSs at T <sub>0</sub> to determinate “low-risk” or “high-risk” for CVD	

Features	
Demographic data	age, gender, race and smoking status
Disease phenotypes	457 phecodes
Drugs by ingredients	866 RxNorm concept codes
Physical measurement	e.g. blood pressure and body mass index (BMI)
Laboratory tests	214 tests
Genetic variants	single nucleotide polymorphisms (SNPs) from catalog
Expression data	6,697 gene expression data

Longitudinal physical measurements and laboratory tests were represented by summarized data (i.e. max, min, median, mean, and counts). Missing data was median-imputed from individuals of the same age, gender, and race. Dummy variables were used to represent missing-imputed data.

FRS was calculated using the most recent measurement, classifying subjects as high or low-risk for CVD:

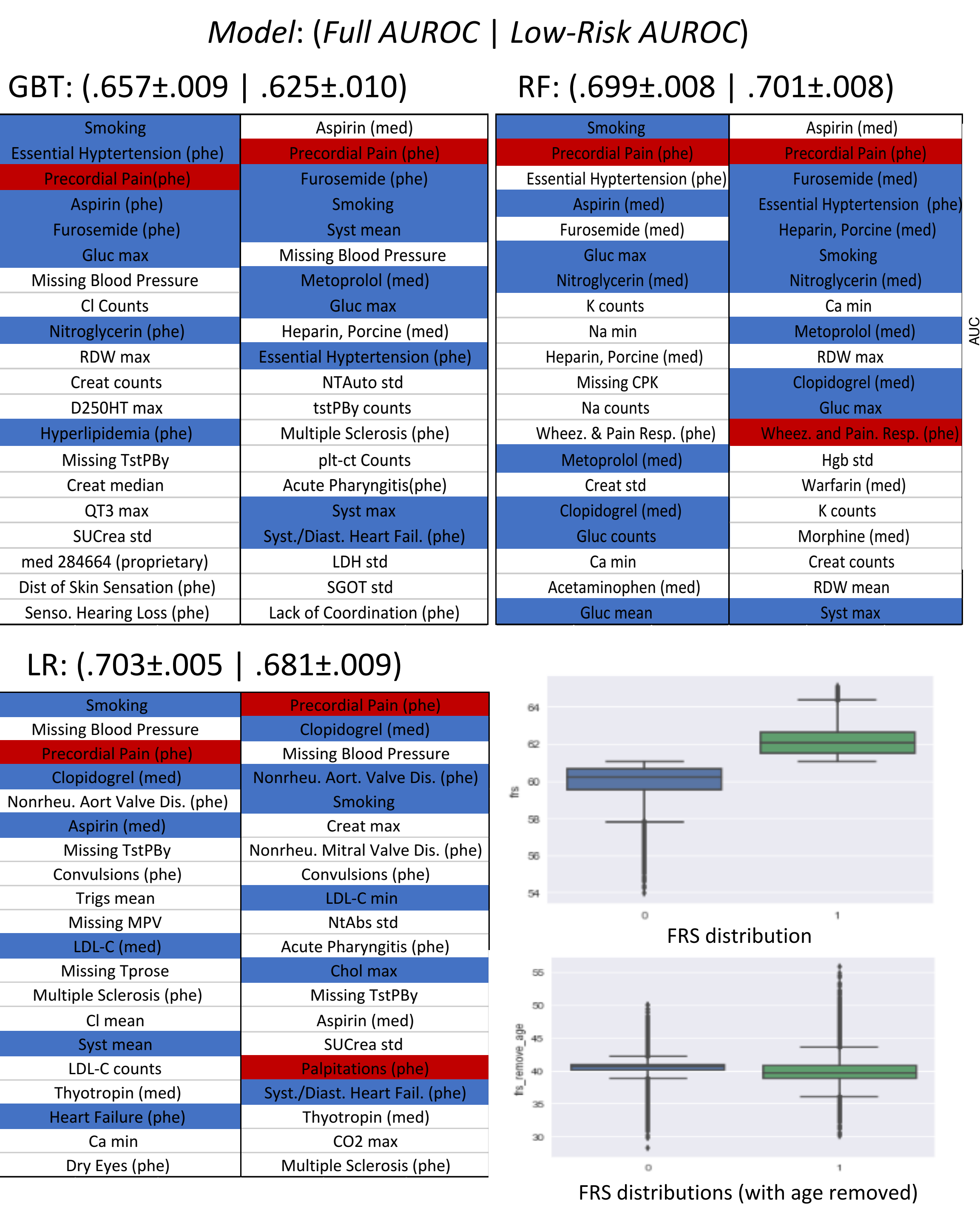
- Full cohort: all cases plus matches controls
- Low-risk cohort: low-risk cases plus matched controls

Three machine learning algorithms implemented and trained for CVD classification:

- Logistic Regression (LR)
- Random Forest (RF)
- Gradient Boosting Trees (GBT)

Classification performance evaluated via 10-fold cross validation AUROC.

**Results**  
Classification performance and the top 20 predictive features for each model are displayed below:



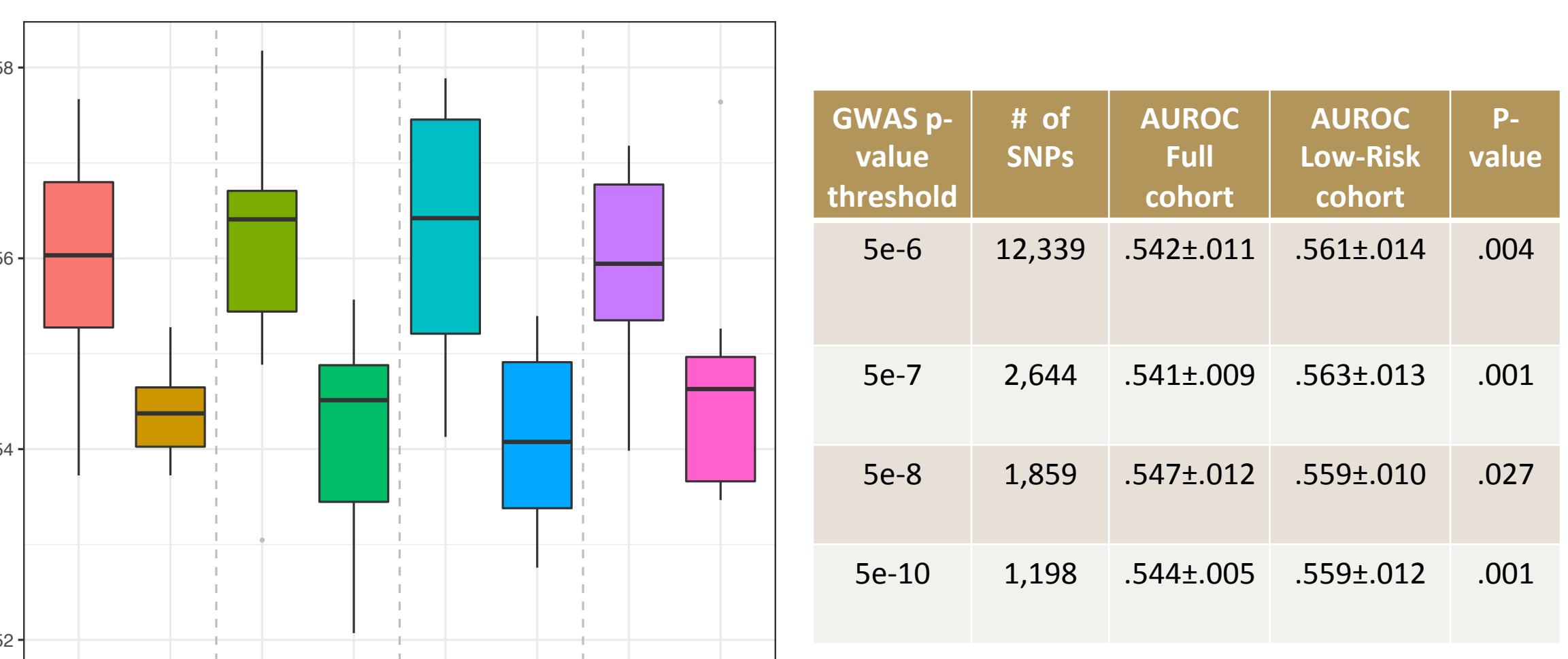
**Key Findings:**  
Similar classification performance and conventional CVD top features when trained on both the full and low-risk cohorts indicate that when age, race, and gender are ignored, the presentation of conventional risk factors is the best predictive indicator for low-risk CVD outcome.

As age is the largest contributor to FRS, our study results indicate that the FRS is highly biased by age, potentially classifying younger individuals presenting conventional symptoms as low-risk despite benefiting from CVD intervention.

**Genetic Factors Play a Bigger Role in CVD Outcome in Low-Risk Individuals Compared to High-Risk**  
We were interested in seeing how the amount of CVD information encoded in SNPs in low-risk individuals compared to high-risk. Individuals were selected from VUMC BioVU, Vanderbilt’s biobank containing fully genotyped individuals:

- 13,900 CVD cases out of 37,573 individuals selected
- 6,553 cases classified as low-risk for CVD by FRS computed with most recent available measurements

Statistically significant SNPs identified from CVD GWAS (Genome-Wide Association Study) with various statistically significant p-value thresholds. CVD classification results from a 3-hidden layer deep neural network trained on SNP data (no demographic information):

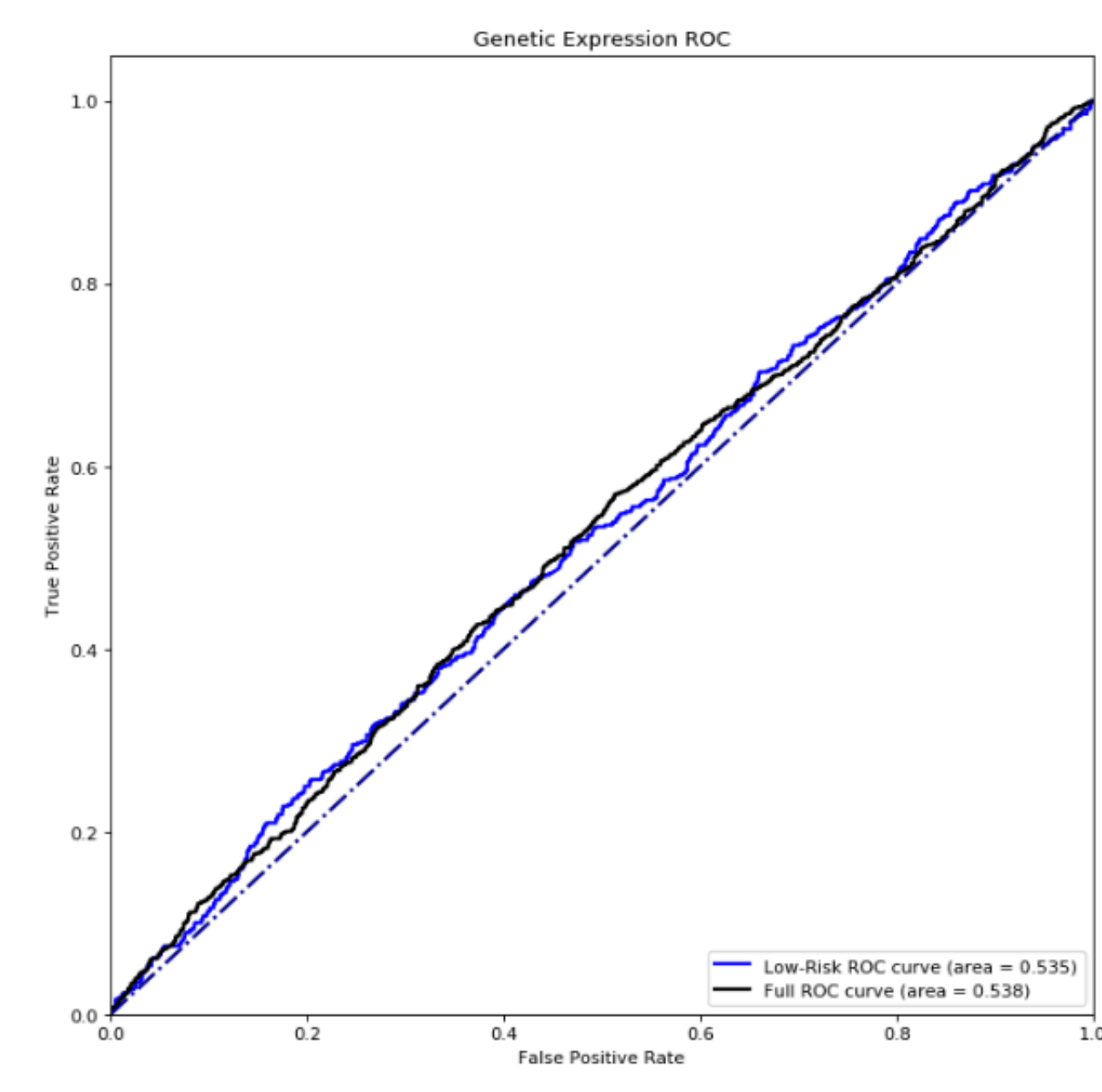


**Key Findings:**  
1) No significant performance gains as more SNPs are included. This means that CVD relevant information is encoded in the 5e-10 GWAS SNPs.  
2) There is a statistically significant performance gain in the prediction on the low-risk cohort over the full cohort pointing to unavoidable genetic causes over long-term lifestyle factors as a possible explanation for why low-risk individuals develop CVD earlier in life or present CVD without notable conventional factors.

**Genetic Expression Data makes no Significant Contribution**  
Synthetic genetic expression data generated from SNP frequency fed into a deep learning neural network

- Representative expression of 6,697 genes from blood
- 12,533 CVD cases (5489 low-risk) out of 26,722 individuals

CVD predictive performance was slightly lower than when trained on raw SNP data. No statistically significant increase in performance for the low-risk cohort over full.



**Conclusions**  
**Longitudinal EHR & Genetic information holds predictive insight for individuals considered low-risk for CVD. Our major conclusions are that the FRS is overly influenced by age causing it to overlook younger individuals with conventional risk factors and that genetic information provides a larger contribution to CVD progression in the low-risk population. Future work includes using stricter definitions for individuals considered low-risk and late-fusion approaches for combining trained predictive models.**

ACKNOWLEDGEMENT: This work was funded by NSF Award 1757644

Thanks for the help from Juan Zhao PhD, Qiping Feng PhD, Josh Denny MD, Patrick Wu BS, and Janey Wang MS