

# Prediction of Disease by Pathway-Based Integrative Genomic and Demographic Analysis

Skanda Koppula<sup>14</sup>, Amin Zollanvari<sup>123</sup>,  
Gil Alterovitz<sup>1234\*</sup>

PRIMES Conference  
May 18, 2013



<sup>1</sup> Center for Biomedical Informatics, Harvard Medical School [Boston, MA 02115].

<sup>2</sup> Children's Hospital Informatics Program at Harvard-MIT Division of Health Science [Boston, MA 02115].

<sup>3</sup> Partners Healthcare Center for Personalized Genetic Medicine [Boston, MA 02115].

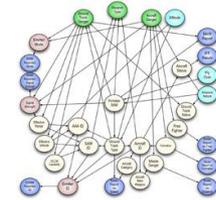
<sup>4</sup> Dept. of Electrical Engineering and Computer Science at MIT [Cambridge, MA 02139].

\* Corresponding author. Contact: [gil@mit.edu](mailto:gil@mit.edu)

# Introduction

## ★ Why prediction-based analysis of data?

- ✓ Flexible model types
- ✓ Gauge effect of feature on phenotype
- ✓ ...effective diagnostic tools!

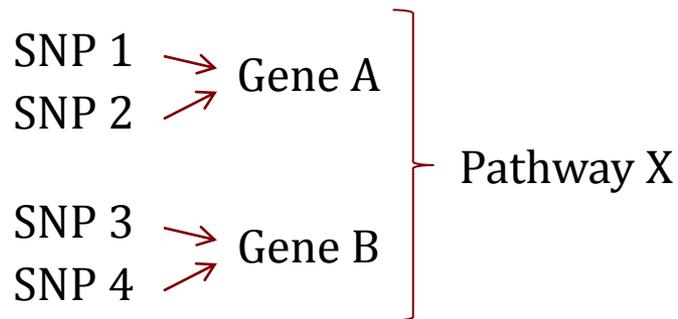


## Introduction

### ★ Why prediction-based analysis of data?

- ✓ Flexible model types
- ✓ Gauge effect of feature on phenotype
- ✓ ...effective diagnostic tools!

### ★ Try analysis on a different level!



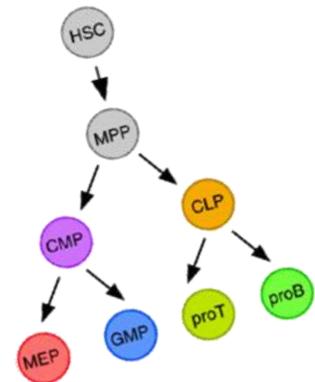
## Introduction

### ☀ Why prediction-based analysis of data?

- ✓ Flexible model types
- ✓ Gauge effect of feature on phenotype
- ✓ ...effective diagnostic tools!

### ☀ Try analysis on a different level?

- ✓ Use inter-gene relations!
- ✓ No black-box around disease mechanism
- ✓ More knowledge about features *with no data*

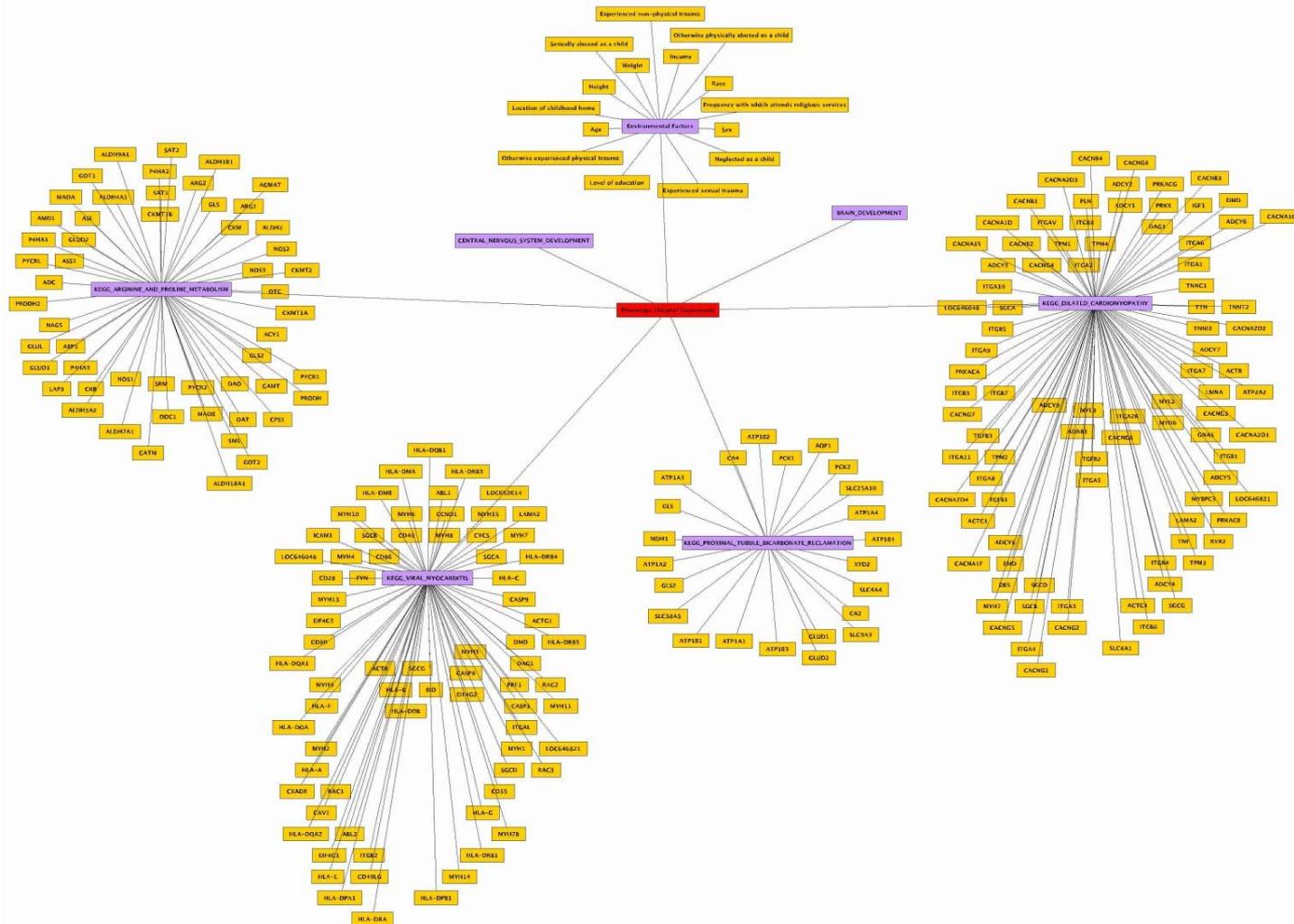


## Introduction

- ✱ Why prediction-based analysis of data?
  - ✓ Flexible models [data type, number of features]
  - ✓ Easy to measure effect of feature on phenotype
  - ✓ Effective diagnostic tool
  
- ✱ Try analysis on a different level?

*Pathway-based predictive models*

# Predictive Framework : TAN and Naïve Bayes



# Alcoholism

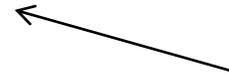
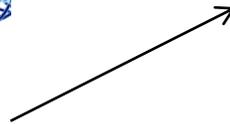
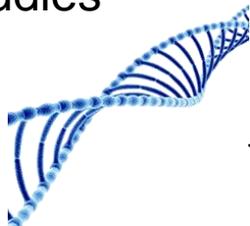
2.5 million

14%

*“increasing consumption of alcohol even in face of adverse consequences”*

---

twin adoption  
studies



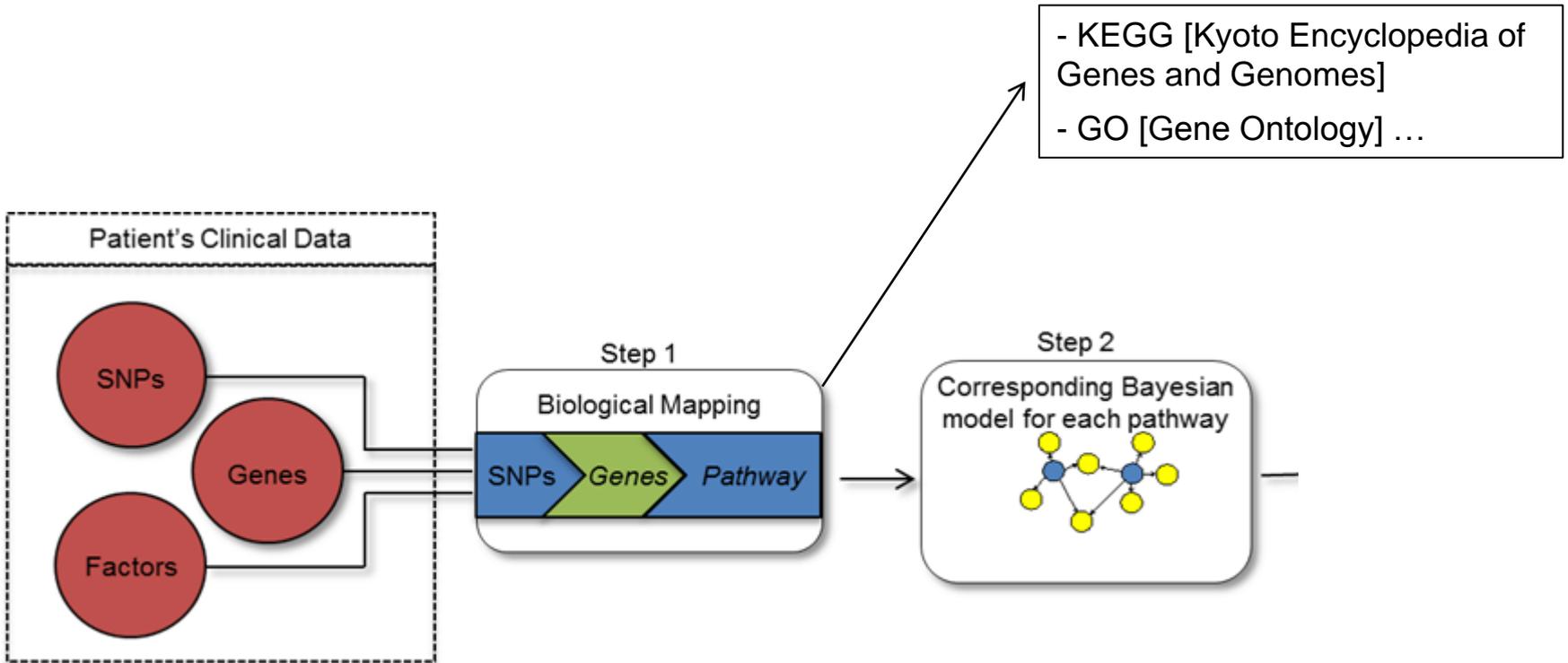
environmental  
studies

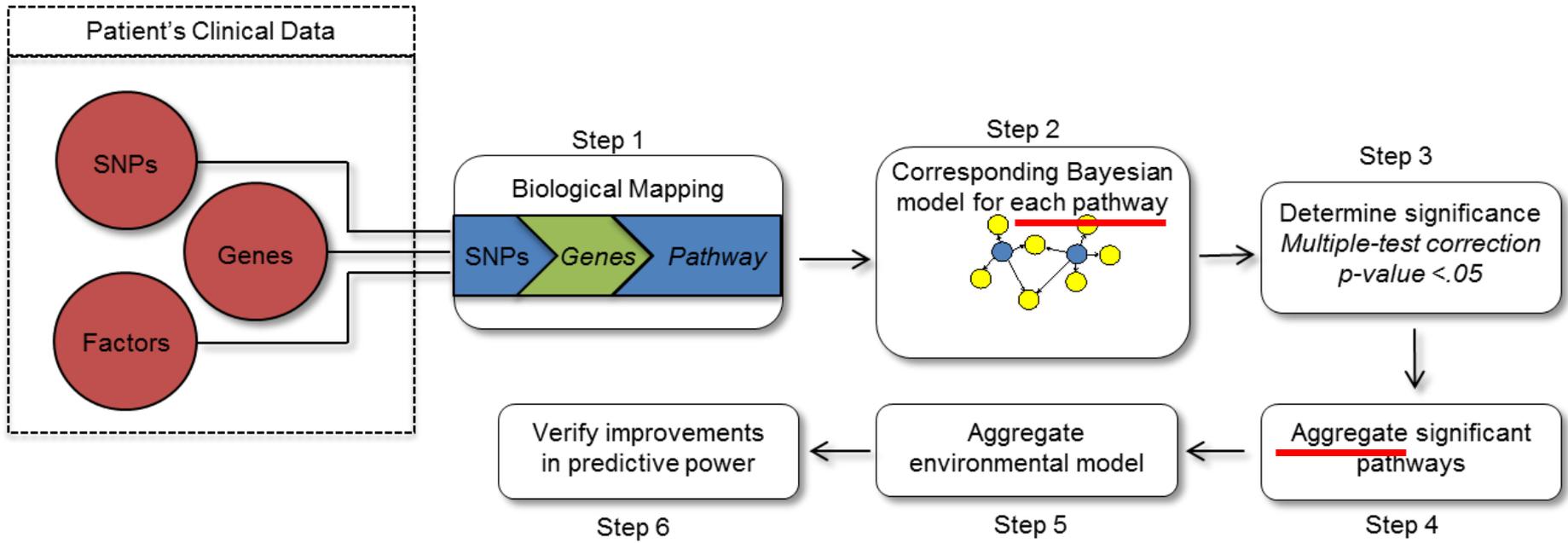


---

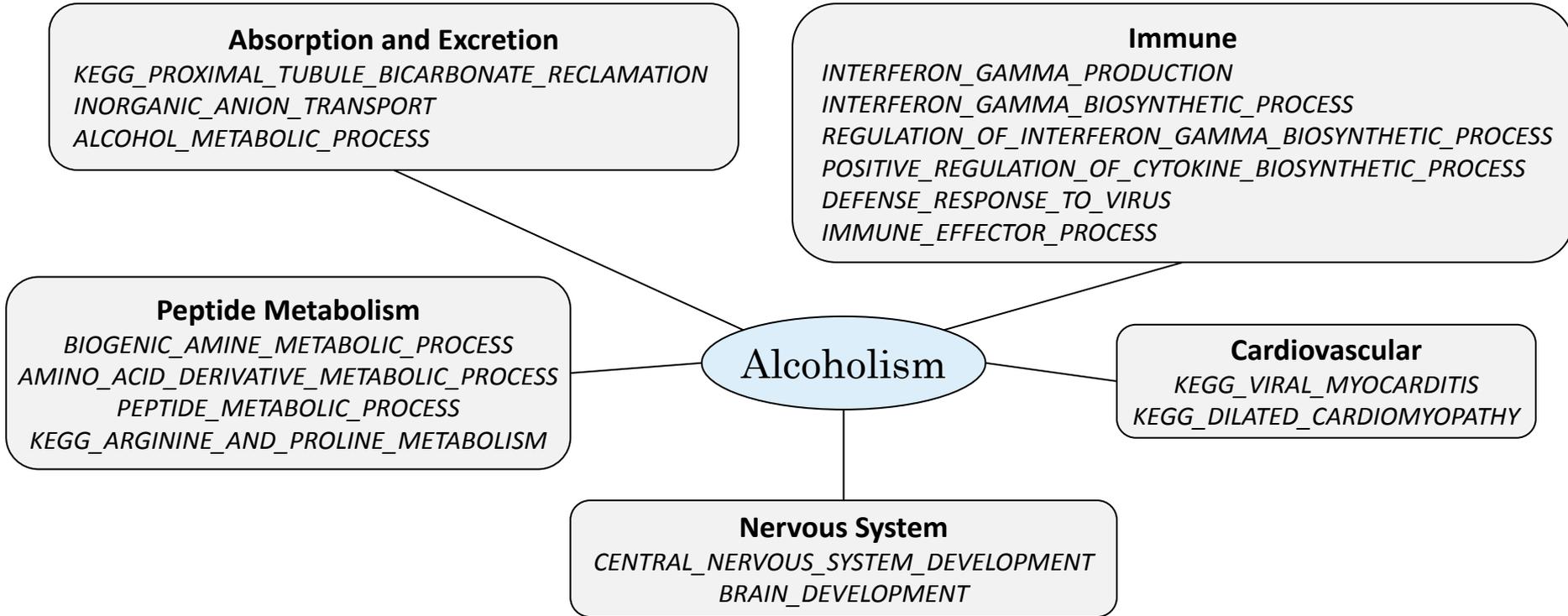
The datasets:

- COGA (1653 patients)
- COGENE (1350 patients)



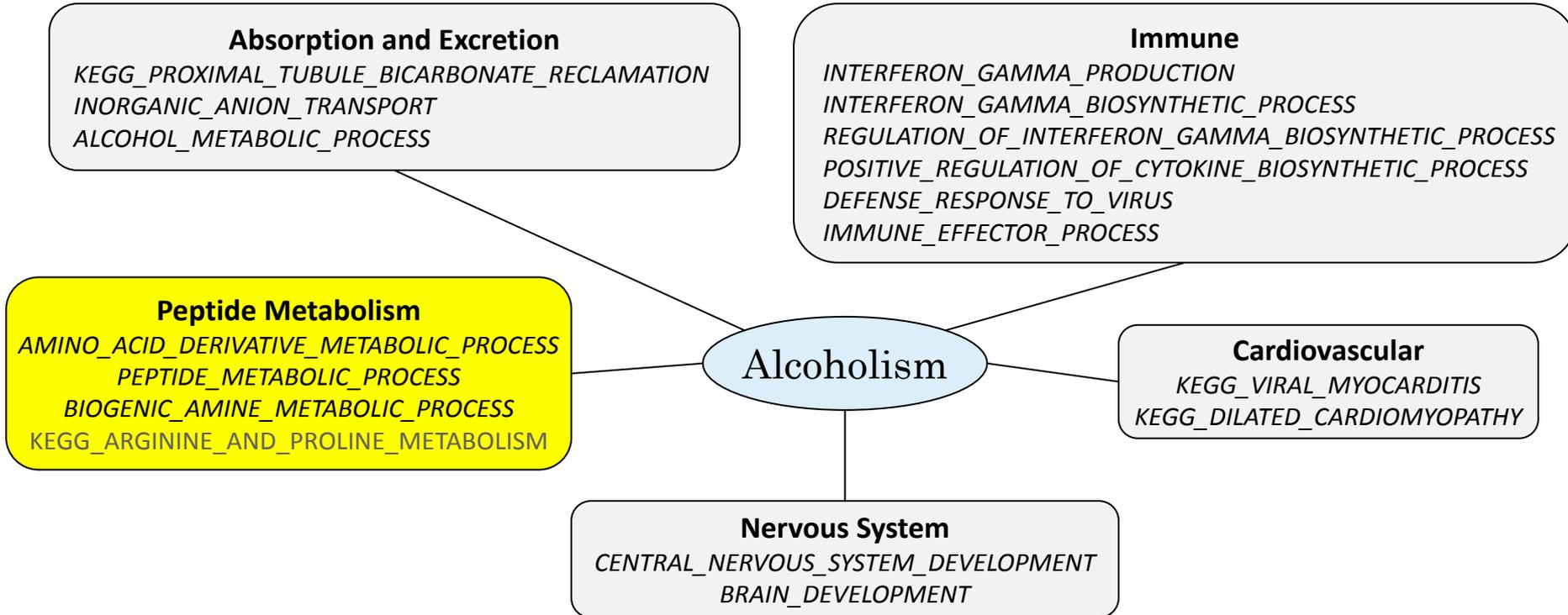


# Genetic-Only Model



AUROC = 0.83  
 $p < 10^{-3}$

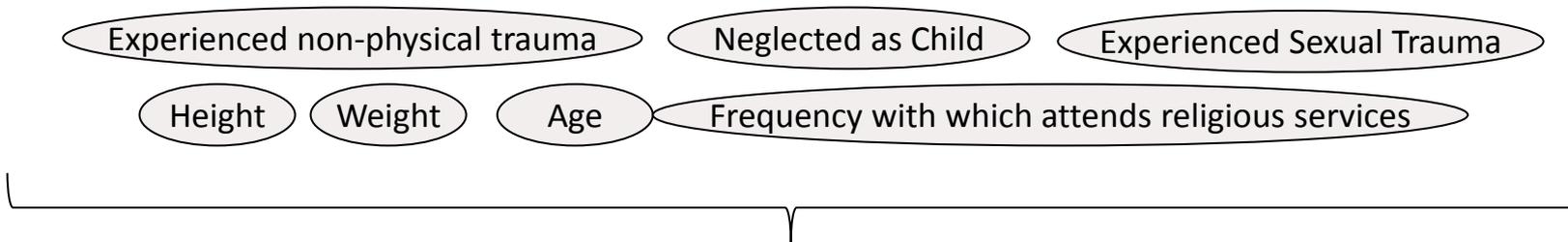
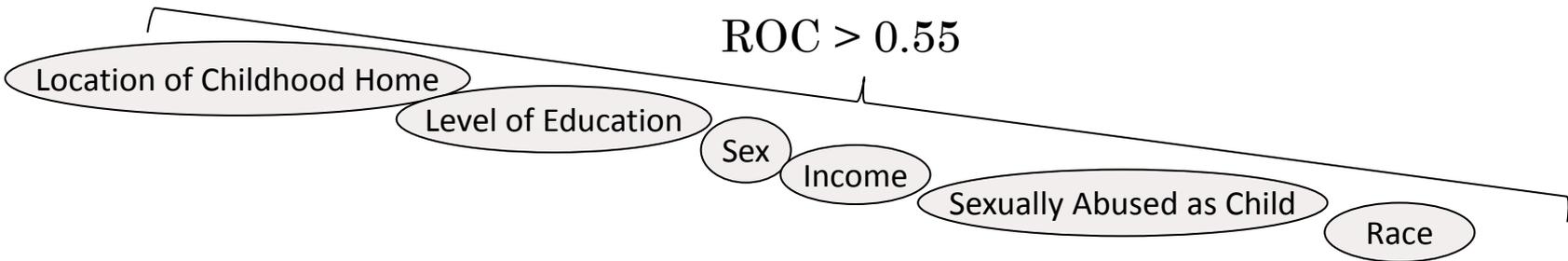
# Genetic-Only Model



 AUROC = 0.83  
 $p < 10^{-3}$

# Genetic-Demographic Model

AUROC = 0.90  
 $p < 10^{-5}$



ROC < 0.55

## Genetic-Demographic Model

- ✱ Increase due to more # features?
  - ✓ No! Replacement increases accuracy by 2.8%
- ✱ Why?
  - ✓ Genes and demo. factors boost each other
    - ☞ *Inorganic Anion Transport contains {CLCNX gene group} on X-chromosome*

## Lung Cancer

Pathway	AUROC
<i>Estrogen receptor regulation (carm1 and -er)</i>	0.75
<i>Eukaryote Translation Initiation Factor (eif4, eif2)</i>	0.73
<i>rnaPathway</i>	0.73
<i>ST_Tumor_Necrosis_Factor_Pathway</i>	0.72
<i>vegfPathway</i>	0.67
<i>MAP00010_Glycolysis_Gluconeogenesis</i>	0.66
<i>P53_UP</i>	0.66

AUROC = 0.85

$p < 10^{-5}$

## Next Steps

1. Insight from inter-feature relationships?
2. Application for layman to use predictive framework?
3. *In vitro* validation of identified pathways
4. Other learning structures?



## Acknowledgements

- ✱ PRIMES program for providing me with this opportunity
  - ✓ Dr. Gerovitch, Professor Etingof, and Professor Khovanova
  
- ✱ Professor Alterovitz
  
- ✱ NIH Grants:
  - ✓ 5R21DA025168-02 (G. Alterovitz)
  - ✓ 1R01HG004836-01 (G. Alterovitz)
  - ✓ 4R00LM009826-03 (G. Alterovitz)

Thank You! Questions?