

# Significance of gene-gene interactions (epistasis)

PSB 2015 Tutorial

Marylyn D Ritchie, PhD

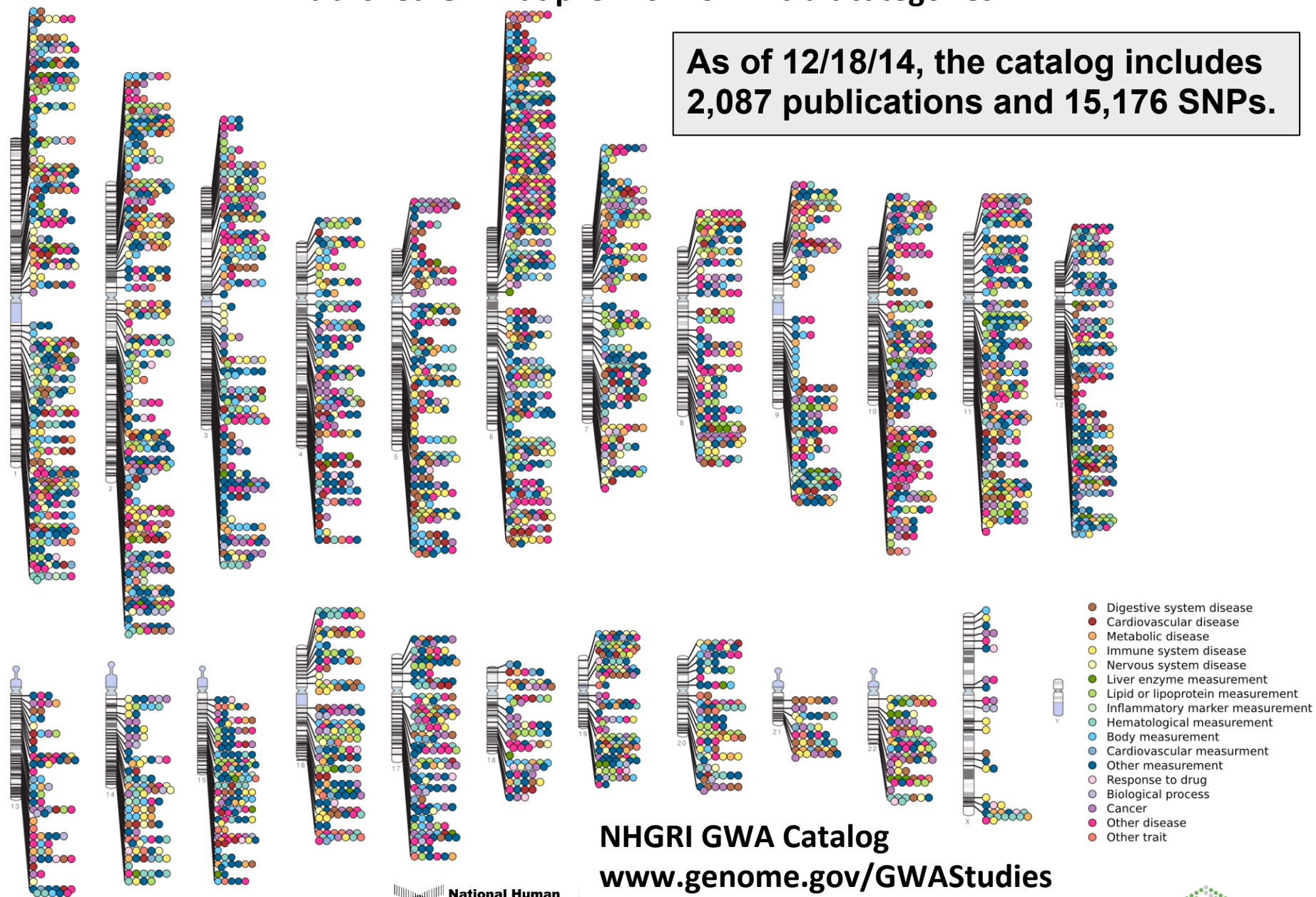
Director, Biomedical and Translational Informatics, Geisinger Clinic

Professor, Biochemistry and Molecular Biology, The Pennsylvania State University

# Published Genome-Wide Associations through 12/2013

Published GWA at  $p \leq 5 \times 10^{-8}$  for 17 trait categories

As of 12/18/14, the catalog includes  
2,087 publications and 15,176 SNPs.

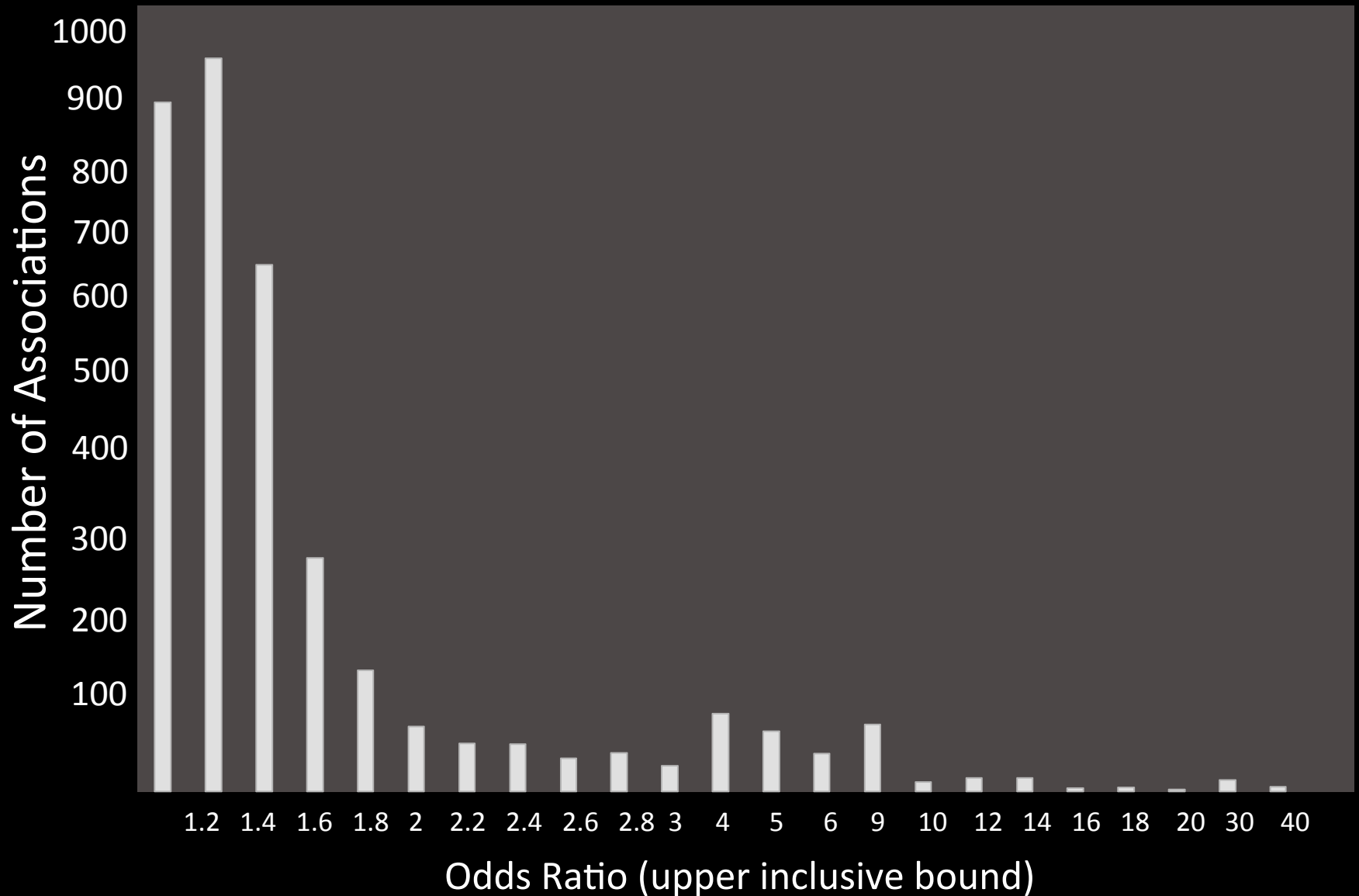


NHGRI GWA Catalog

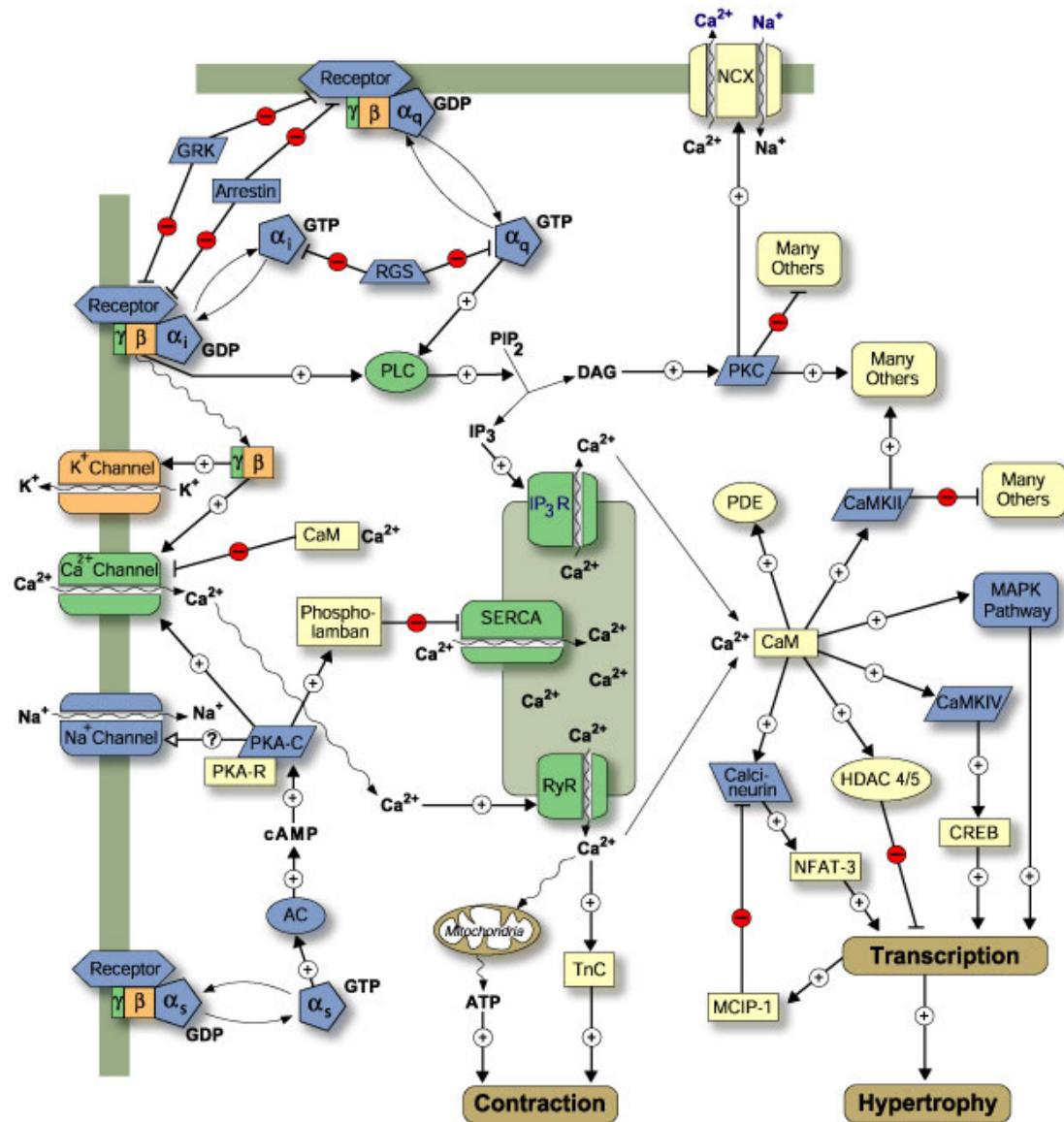
[www.genome.gov/GWAStudies](http://www.genome.gov/GWAStudies)

[www.ebi.ac.uk/fgpt/gwas/](http://www.ebi.ac.uk/fgpt/gwas/)

# Distribution of Effects



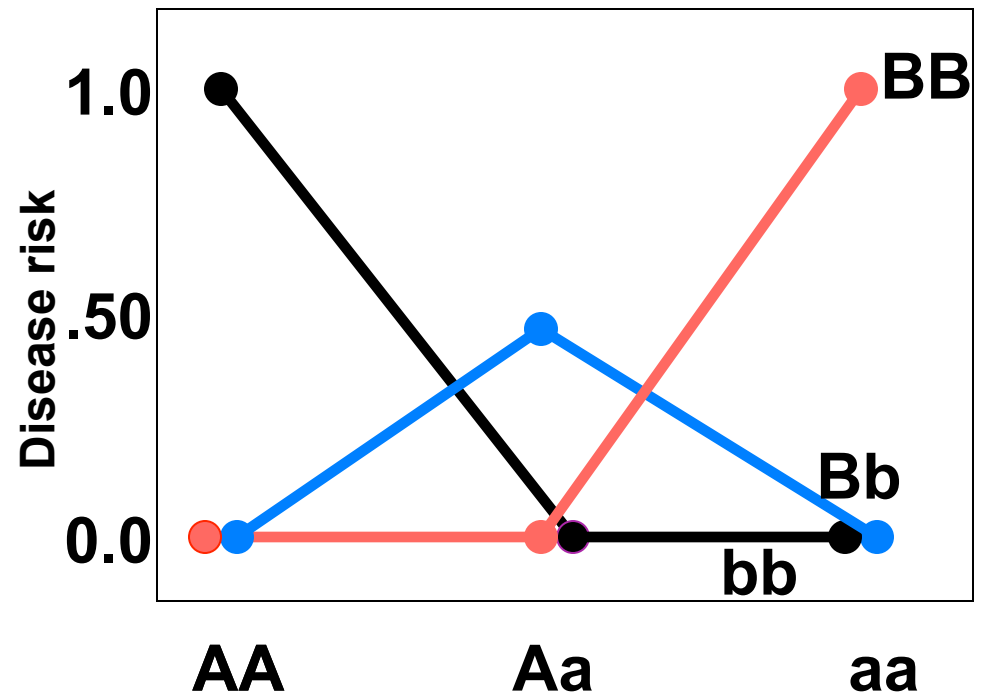
# Biology is complex



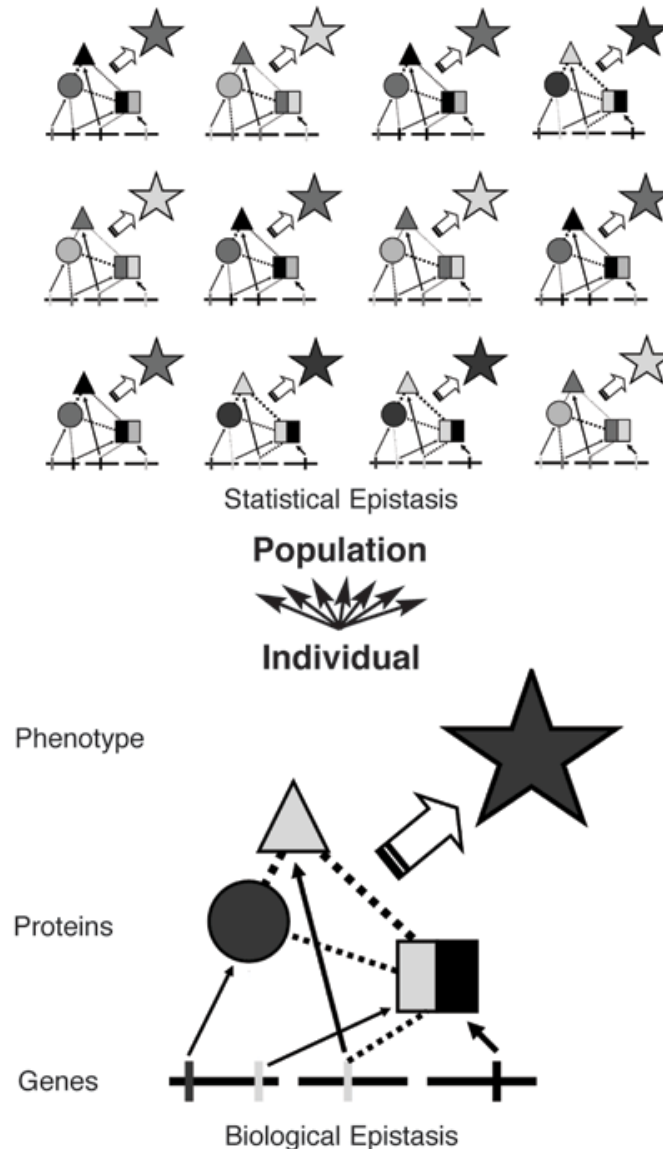
# Epistasis

- Epistasis – two or more genes interacting in a non-additive manner to confer disease risk; gene-gene interactions

Genotype	p(D)
AABB	0.0
AABb	0.0
AAbb	1.0
AaBB	0.0
AaBb	.50
Aabb	0.0
aaBB	1.0
aaBb	0.0
aabb	0.0

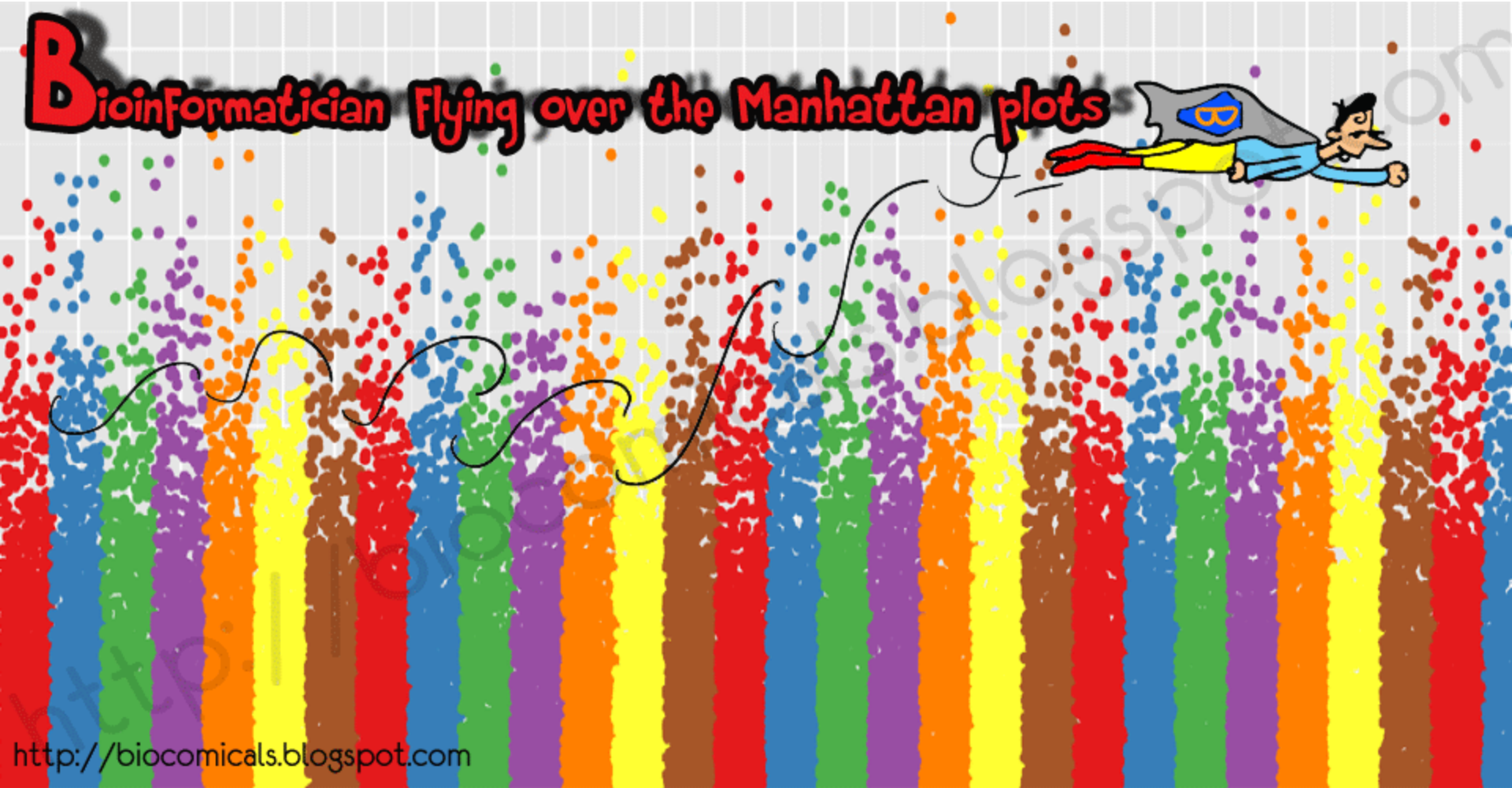


# Statistical Epistasis vs. Biological Epistasis





# Traditional Approach



# Traditional Statistical Approaches

## Genetic Epidemiology - Association Analysis

- Typically one marker or SNP at a time to detect loci exhibiting main effects
- Follow-up with an analysis to detect interactions between the main effect loci
- Some studies attempt to detect pair-wise interactions even without main effects
- Higher dimensions are usually not possible with traditional methods



# Traditional Statistical Approaches

## Genetic Epidemiology - Association Analysis

### ■ Logistic Regression

- ◆ Small sample size can result in biased estimates of regression coefficients and can result in spurious associations (Concato et al. 1993)
- ◆ Need at least 10 cases or controls per independent variable to have enough statistical power (Peduzzi et al. 1996)
- ◆ Curse of dimensionality is the problem (Bellman 1961)

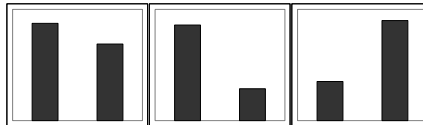
# Curse of Dimensionality

**N = 100**

**50 Cases, 50 Controls**

**SNP 1**

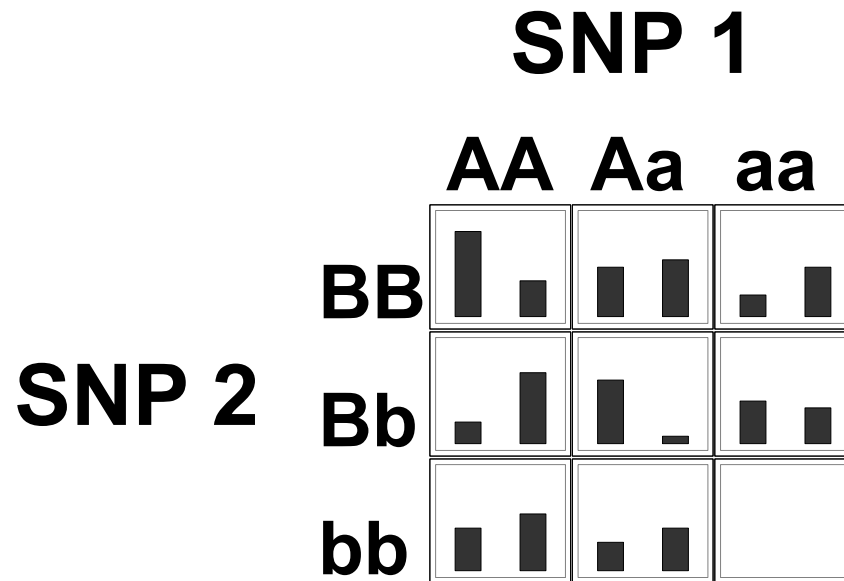
**AA Aa aa**



# Curse of Dimensionality

N = 100

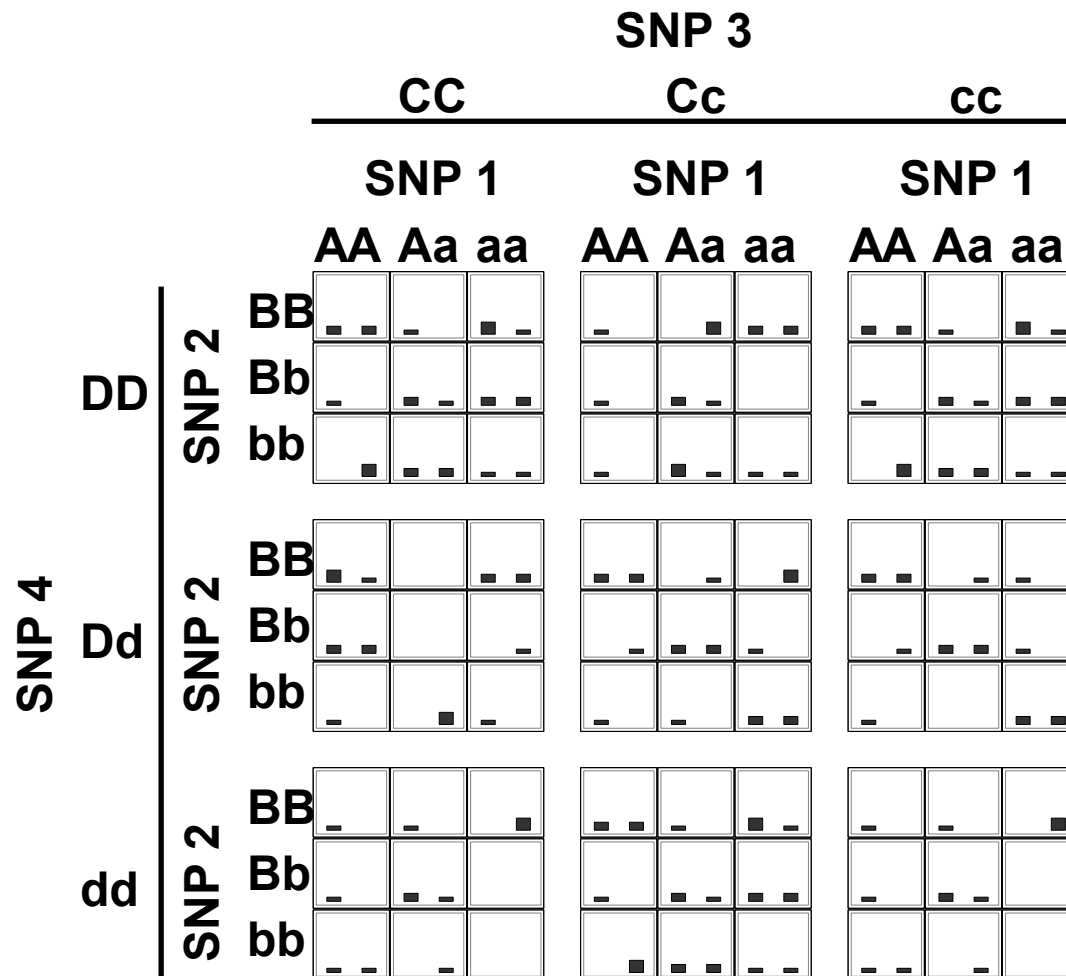
50 Cases, 50 Controls



# Curse of Dimensionality

N = 100

50 Cases, 50 Controls

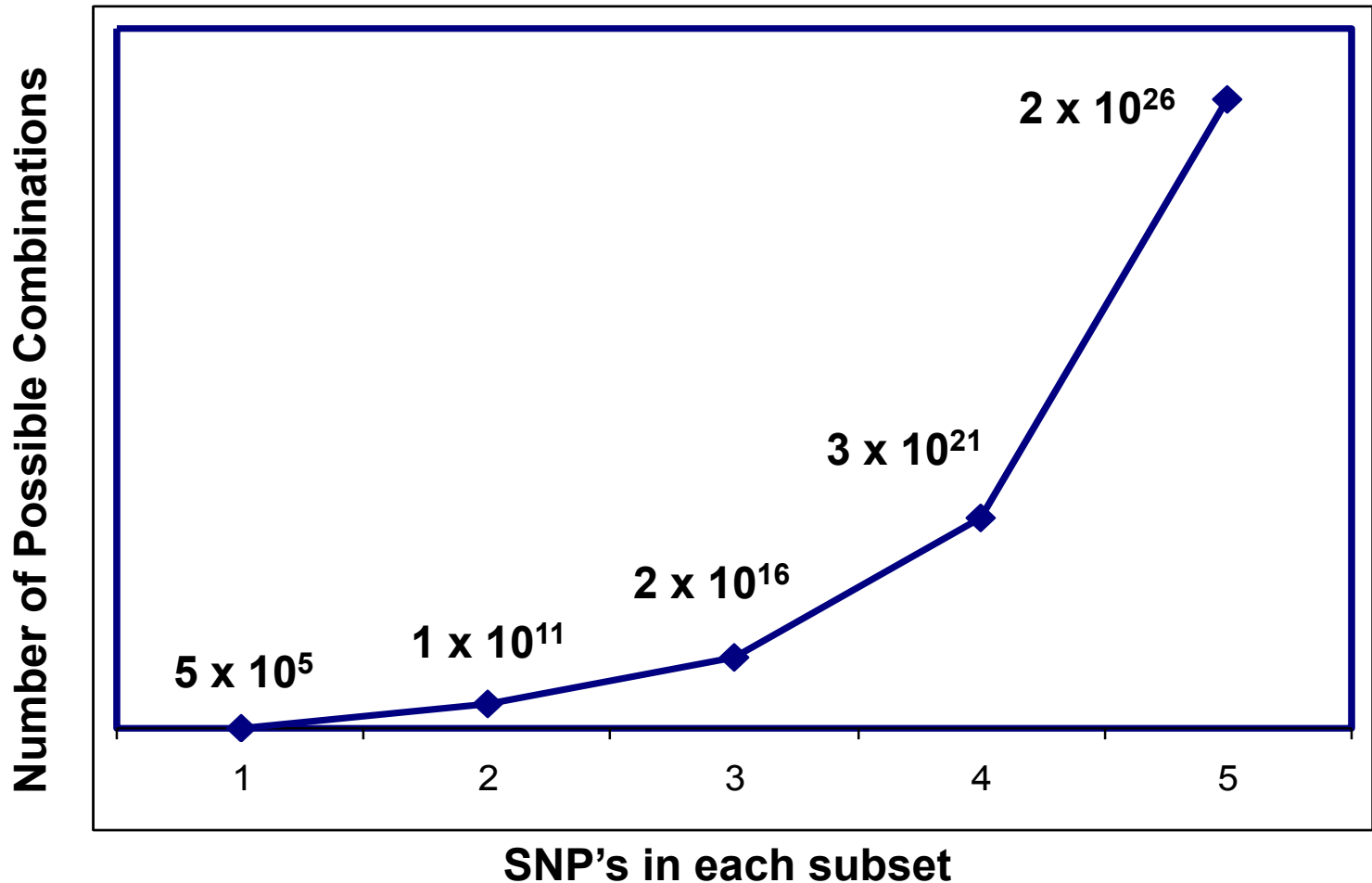


**If interactions with minimal main effects are the norm rather than the exception, can we analyze all possible combinations of loci with traditional approaches to detect purely interaction effects ?**

**NO**

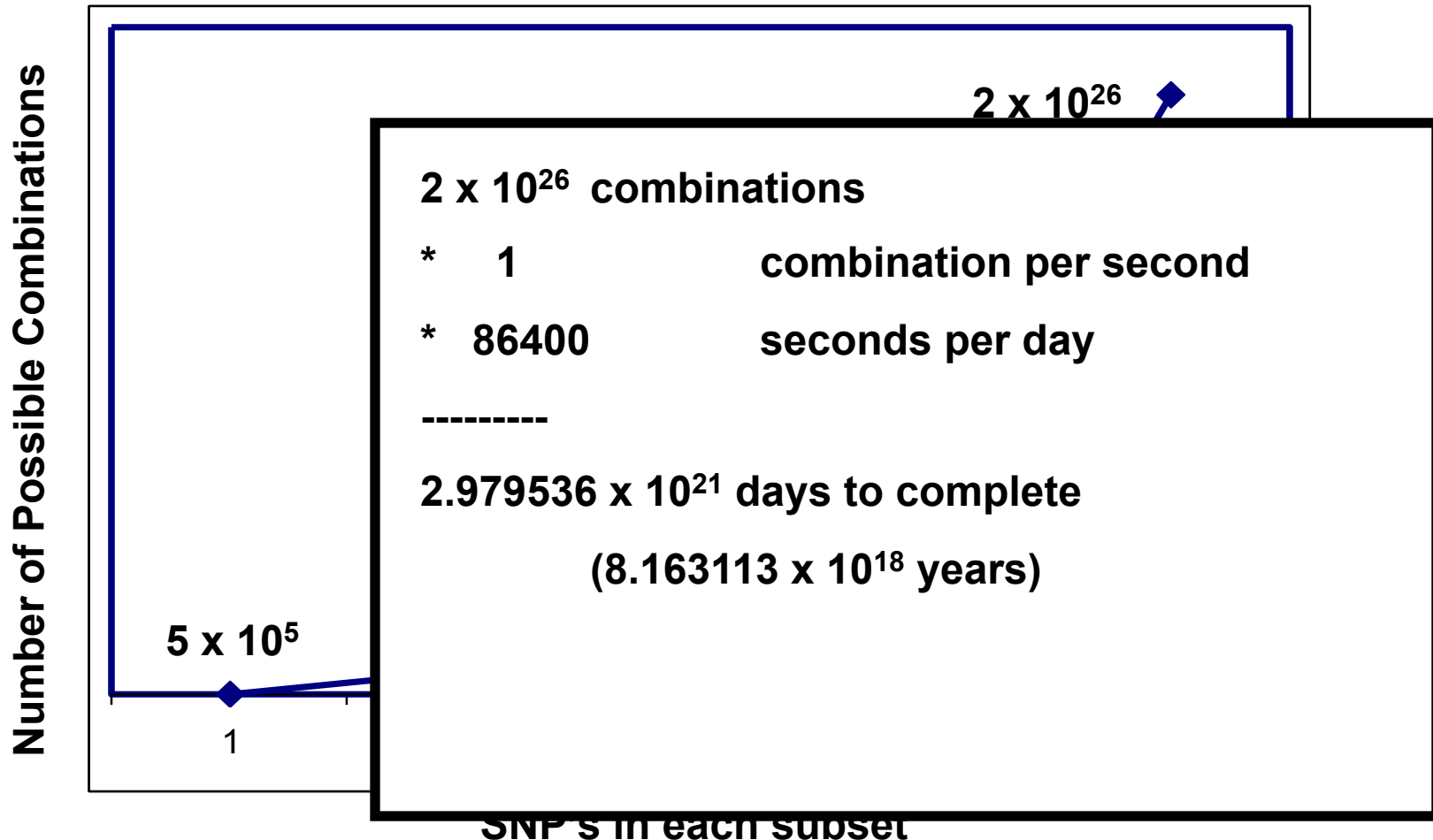
# How many combinations are there?

- ~500,000 SNPs to span the genome (HapMap)



# How many combinations are there?

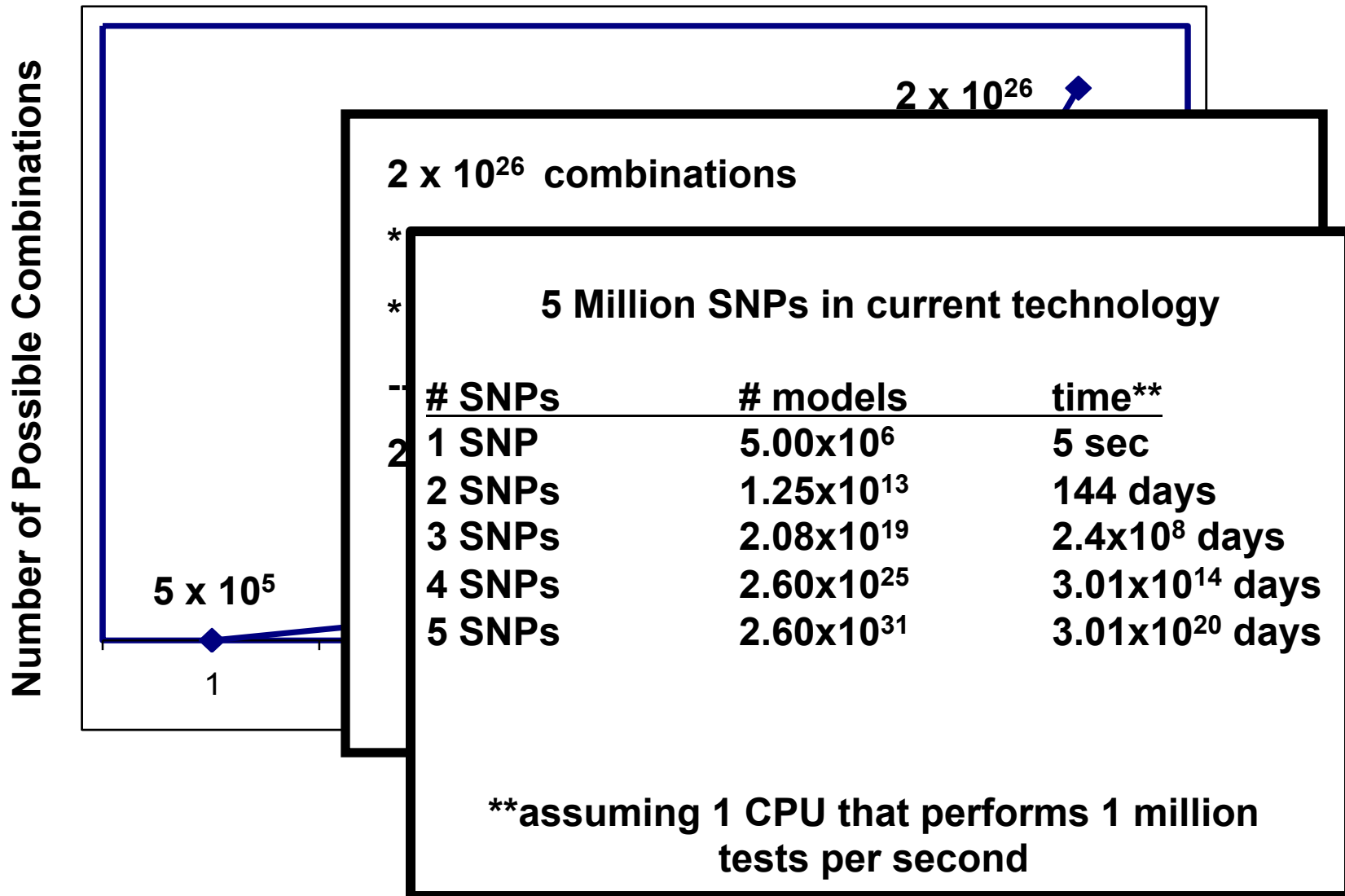
- ~500,000 SNPs to span the genome (HapMap)





# How many combinations are there?

- ~500,000 SNPs to span the genome (HapMap)



# THE BIG BANG THEORY

**$5.47 \times 10^{12}$  days**

TIME  
BEGINS

ONE  
SECOND

PRESENT  
DAY

Time	$10^{-43}$ sec.	$10^{-32}$ sec.	$10^{-6}$ sec.	3 min.	300,000 yrs.	1 billion yrs.	15 billion yrs.
Temperature		$10^{27}^{\circ}\text{C}$	$10^{13}^{\circ}\text{C}$	$10^8^{\circ}\text{C}$	$10,000^{\circ}\text{C}$	$-200^{\circ}\text{C}$	$-270^{\circ}\text{C}$

**1** The cosmos goes through a superfast "inflation," expanding from the size of an atom to that of a grapefruit in a tiny fraction of a second

**2** Post-inflation, the universe is a seething, hot soup of electrons, quarks and other particles

**3** A rapidly cooling cosmos permits quarks to clump into protons and neutrons

**4** Still too hot to form into atoms, charged electrons and protons prevent light from shining; the universe is a superhot fog

**5** Electrons combine with protons and neutrons to form atoms, mostly hydrogen and helium. Light can finally shine

**6** Gravity makes hydrogen and helium gas coalesce to form the giant clouds that will become galaxies; smaller clumps of gas collapse to form the first stars

**7** As galaxies cluster together under gravity, the first stars die and spew heavy elements into space; these will eventually form into new stars and planets

NOTE: The numbers in cosmology are so great and the numbers in subatomic physics are so small that it is often necessary to express them in exponential form. Ten multiplied by itself, or 100, is written as  $10^2$ . One thousand is written as  $10^3$ . Similarly, one-tenth is  $10^{-1}$ , and one-hundredth is  $10^{-2}$ .

Source: The Birth of the Universe; The Kingfisher Young People's Book of Space

TIME Graphic by Ed Gabel

# Traditional Approach

## ■ Advantages

- ◆ Computationally feasible
- ◆ Easy to interpret

## ■ Disadvantages

- ◆ Genes must have large main effects
- ◆ Difficult to detect genes if interactions with other genetic and environmental factors are important
- ◆ CANNOT do an exhaustive search

# New Statistical Approaches

- Review paper

For reprint orders, please contact:  
reprints@futuremedicine.com



## Novel methods for detecting epistasis in pharmacogenomics studies

Alison A Motsinger<sup>1</sup>,  
Marylyn D Ritchie<sup>2</sup> &  
David M Reif<sup>†</sup>

<sup>†</sup>Author for correspondence

<sup>1</sup>North Carolina State

University,

Bioinformatics Research

Center,

Department of Statistics,

Raleigh,

NC 27695, USA

<sup>2</sup>Vanderbilt University,

Center for Human Genetics

Research,

Department of Molecular

The importance of gene–gene and gene–environment interactions in the underlying genetic architecture of common, complex phenotypes is gaining wide recognition in the field of pharmacogenomics. In epidemiological approaches to mapping genetic variants that predict drug response, it is important that researchers investigate potential epistatic interactions. In the current review, we discuss data-mining tools available in genetic epidemiology to detect such interactions and appropriate applications. We survey several classes of novel methods available and present an organized collection of successful applications in the literature. Finally, we provide guidance as to how to incorporate these novel methods into a genetic analysis. The overall goal of this paper is to aid researchers in developing an analysis plan that accounts for gene–gene and gene–environment in their own work.

- Pharmacogenomics. 2007 8(9) :1229-41.
- Reviews approximately 40 methods developed to detect gene-gene and gene-environment interactions

# New Statistical Approaches

Chen *et al. BMC Genomics* 2011, **12**:344  
<http://www.biomedcentral.com/1471-2164/12/344>



## METHODOLOGY ARTICLE

## Open Access

# Comparative analysis of methods for detecting interacting loci

Li Chen<sup>1</sup>, Guoqiang Yu<sup>1</sup>, Carl D Langefeld<sup>2</sup>, David J Miller<sup>3</sup>, Richard T Guy<sup>2</sup>, Jayaram Raghuram<sup>3</sup>, Xiguo Yuan<sup>1</sup>, David M Herrington<sup>4</sup> and Yue Wang<sup>1\*</sup>

### Abstract

**Background:** Interactions among genetic loci are believed to play an important role in disease risk. While many methods have been proposed for detecting such interactions, their relative performance remains largely unclear, mainly because different data sources, detection performance criteria, and experimental protocols were used in the papers introducing these methods and in subsequent studies. Moreover, there have been very few studies strictly focused on comparison of existing methods. Given the importance of detecting gene-gene and gene-environment interactions, a rigorous, comprehensive comparison of performance and limitations of available interaction detection methods is warranted.

# New Statistical Approaches

Shang et al. *BMC Bioinformatics* 2011, **12**:475  
<http://www.biomedcentral.com/1471-2105/12/475>



## METHODOLOGY ARTICLE

## Open Access

# Performance analysis of novel methods for detecting epistasis

Junliang Shang<sup>1\*</sup>, Junying Zhang<sup>1\*</sup>, Yan Sun<sup>2</sup>, Dan Liu<sup>1</sup>, Daojun Ye<sup>1</sup> and Yaling Yin<sup>1,3</sup>

### Abstract

**Background:** Epistasis is recognized fundamentally important for understanding the mechanism of disease-causing genetic variation. Though many novel methods for detecting epistasis have been proposed, few studies focus on their comparison. Undertaking a comprehensive comparison study is an urgent task and a pathway of the methods to real applications.

**Results:** This paper aims at a comparison study of epistasis detection methods through applying related software packages on datasets. For this purpose, we categorize methods according to their search strategies, and select five representative methods (TEAM, BOOST, SNPRuler, AntEpiSeeker and epiMODE) originating from different underlying techniques for comparison. The methods are tested on simulated datasets with different size, various epistasis



# Simple Fitness Landscape

**Fitness**



Mt. Fuji

**Model**



# Complex Fitness Landscape

Waimea Canyon

**Fitness**



**Model**

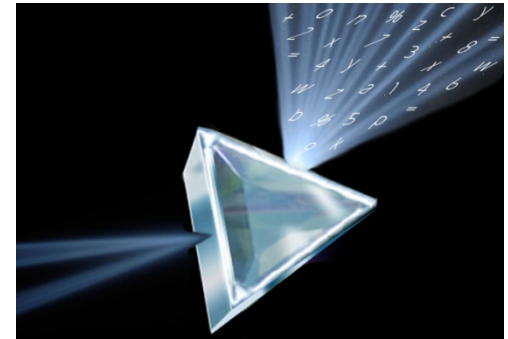
# Epistasis in GWAS Data

- ~~■ Exhaustive evaluation~~
- Evaluate interactions in top hits from single-SNP analysis
- Use prior biological knowledge to evaluate specific combinations – “Candidate Epistasis”

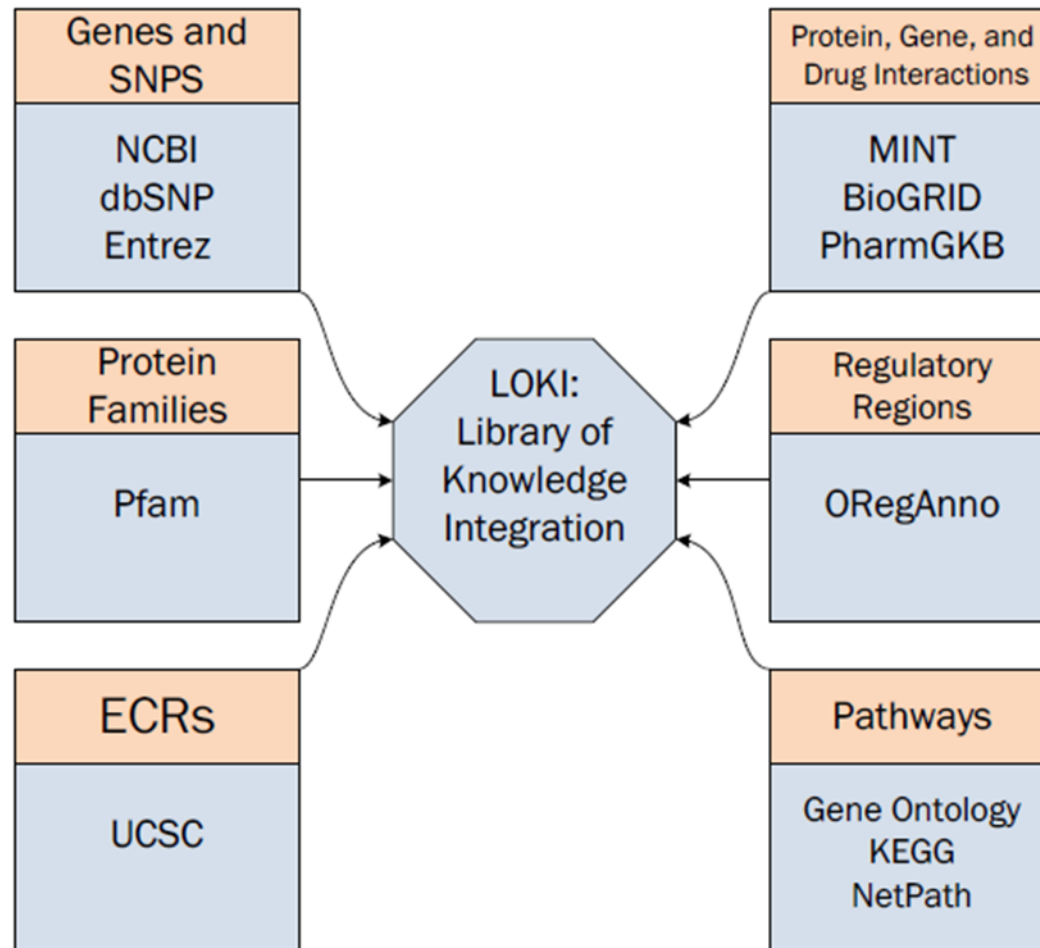
Goal: to build biologically plausible models of gene-gene interactions to test for association using an automated bioinformatics tool based on biological features

# The Biofilter

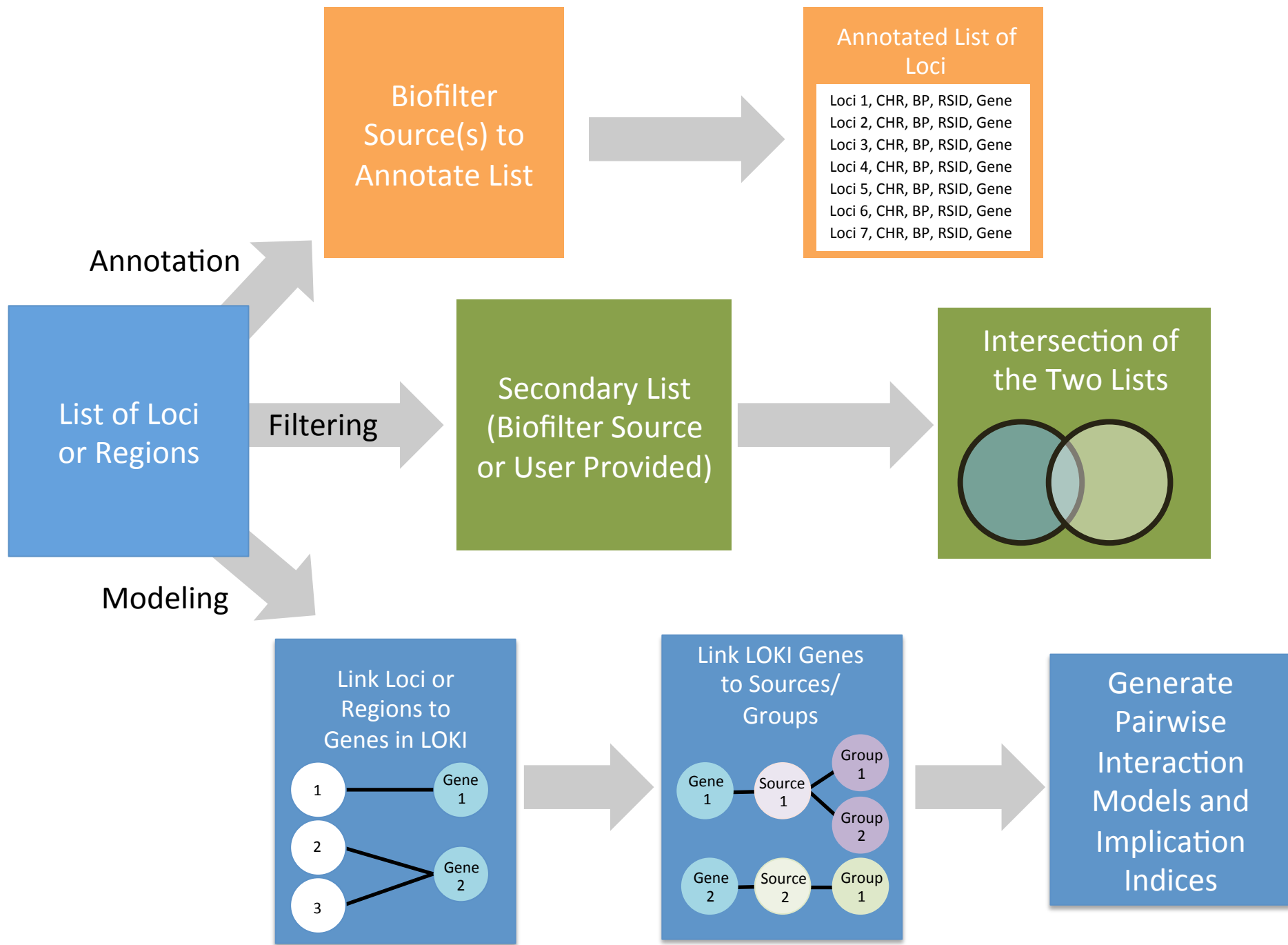
- Use publicly available databases to establish relationships between gene-products
- Suggestions of biological epistasis between genes
- Integrating information from the genome, transcriptome, and proteome into analysis



# LOKI: Library of Knowledge Integration



Bush WS, Dudek SM, Ritchie MD. Biofilter: a knowledge-integration system for the multi-locus analysis of genome-wide association studies. *Pacific Symposium on Biocomputing*, 368-79 (2009).



# Summary

- Gene-gene interactions are important components of complex trait genetic architecture
- Gene-gene interactions are challenging to detect:
  - Due to data sparseness in high dimensions
  - Due to the combinatorics of the search
  - Due to complexity
- Much research is ongoing to develop novel methods and strategies to address these issues