# Data Tsunami as a Limiting Step in Using the All Omics Approach

## ESMO Asia 2015

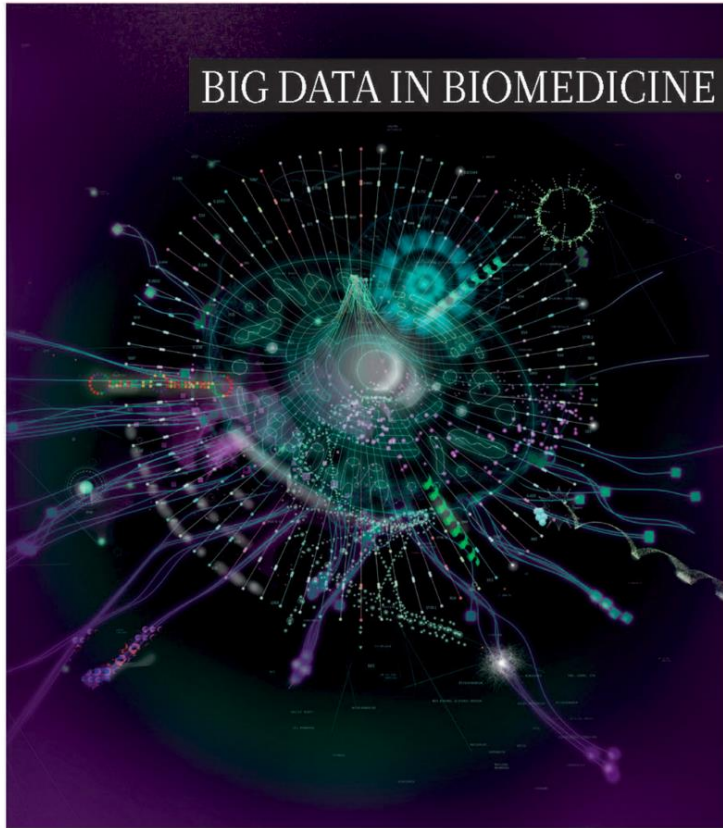**Yu Shyr, Ph.D.**

**Vanderbilt University**

**yu.shyr@vanderbilt.edu**

Advances in sequencing technology have triggered a **tsunami of genomic data**, and these are joined by waves of information from other '-omics' studies, clinical trials and patient records. **Analysis of this big data is launching the era of precision medicine** — **but enormous scientific, engineering and institutional challenges remain**.

**Original Investigation**

# Using Multiplexed Assays of Oncogenic Drivers in Lung Cancers to Select Targeted Drugs

Mark G. Kris, MD , et. al.



A | Patients with an oncogenic driver mutation who did and did not receive targeted therapy, and patients without an ocogenic driver

# Highlights

- **Overview of the BIG data in biomedical research**

  - **Omics data**

  - **EHR data**

  - **Data from patients & other sources**

- **Analytical challenges & tasks**

- **Future of the BIG data in biomedical research**

# Omics biomedical research

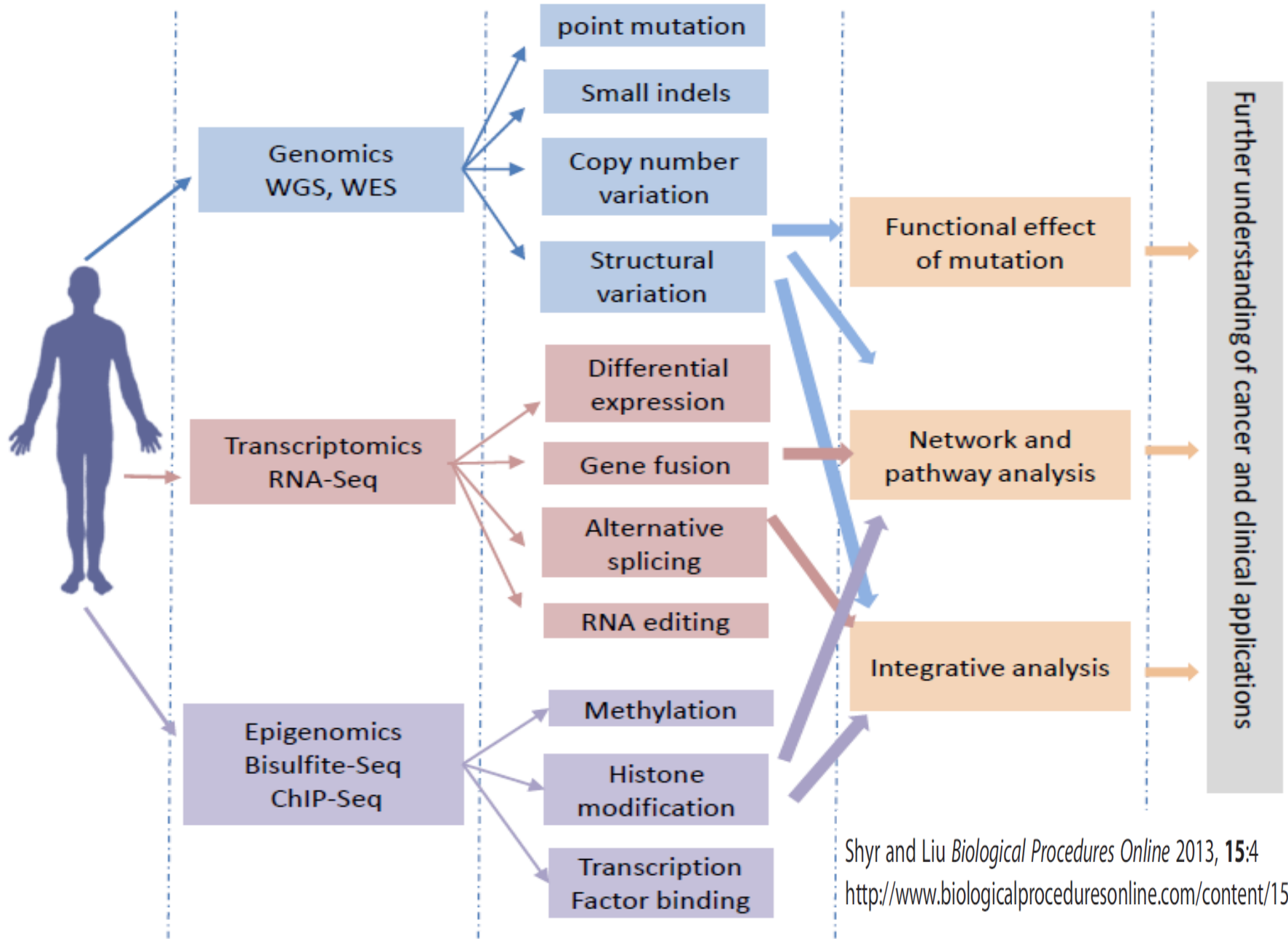- **Microarray: cDNA (about <span style="color:blue">5,000</span> variables), Affymetrix U133 Plus 2.0 (about <span style="color:blue">45,000</span> variables)**

- **SNPs (about <span style="color:blue">500,000 – 2,000,000</span> variables)**
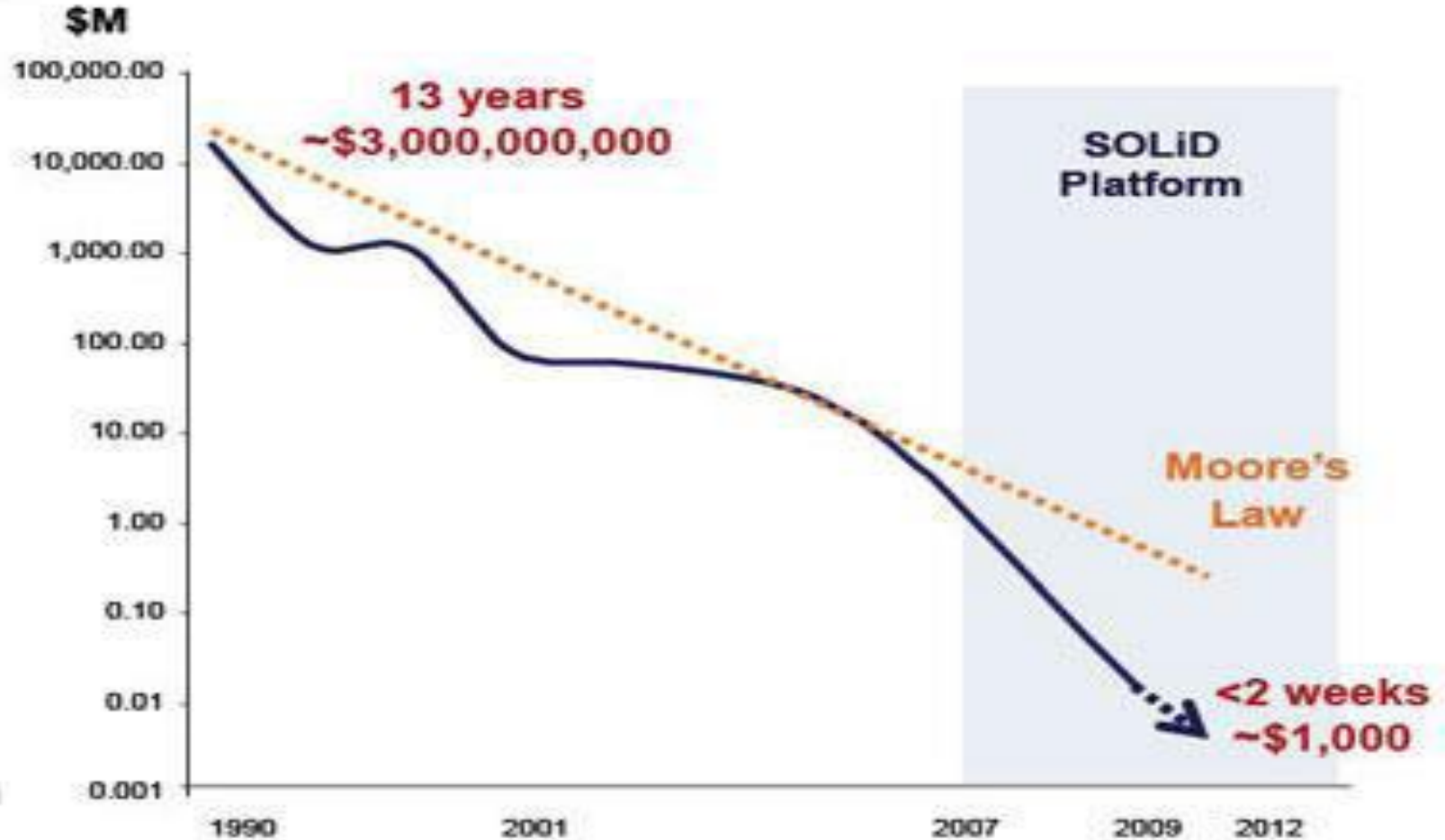
- **Next Generation Sequencing (?)**

Patient | Technologies | Data Analysis | Integration and interpretation

**Patient**

**Technologies**
- Genomics WGS, WES
- Transcriptomics RNA-Seq
- Epigenomics Bisulfite-Seq ChIP-Seq

**Data Analysis**
- point mutation
- Small indels
- Copy number variation
- Structural variation
- Differential expression
- Gene fusion
- Alternative splicing
- RNA editing
- Methylation
- Histone modification
- Transcription Factor binding

**Integration and interpretation**
- Functional effect of mutation
- Network and pathway analysis
- Integrative analysis
- Further understanding of cancer and clinical applications

Shyr and Liu *Biological Procedures Online* 2013, **15**:4
http://www.biologicalproceduresonline.com/content/15/1/4

# Storage of the Data?

- **cDNA, Microarray, SNPs**
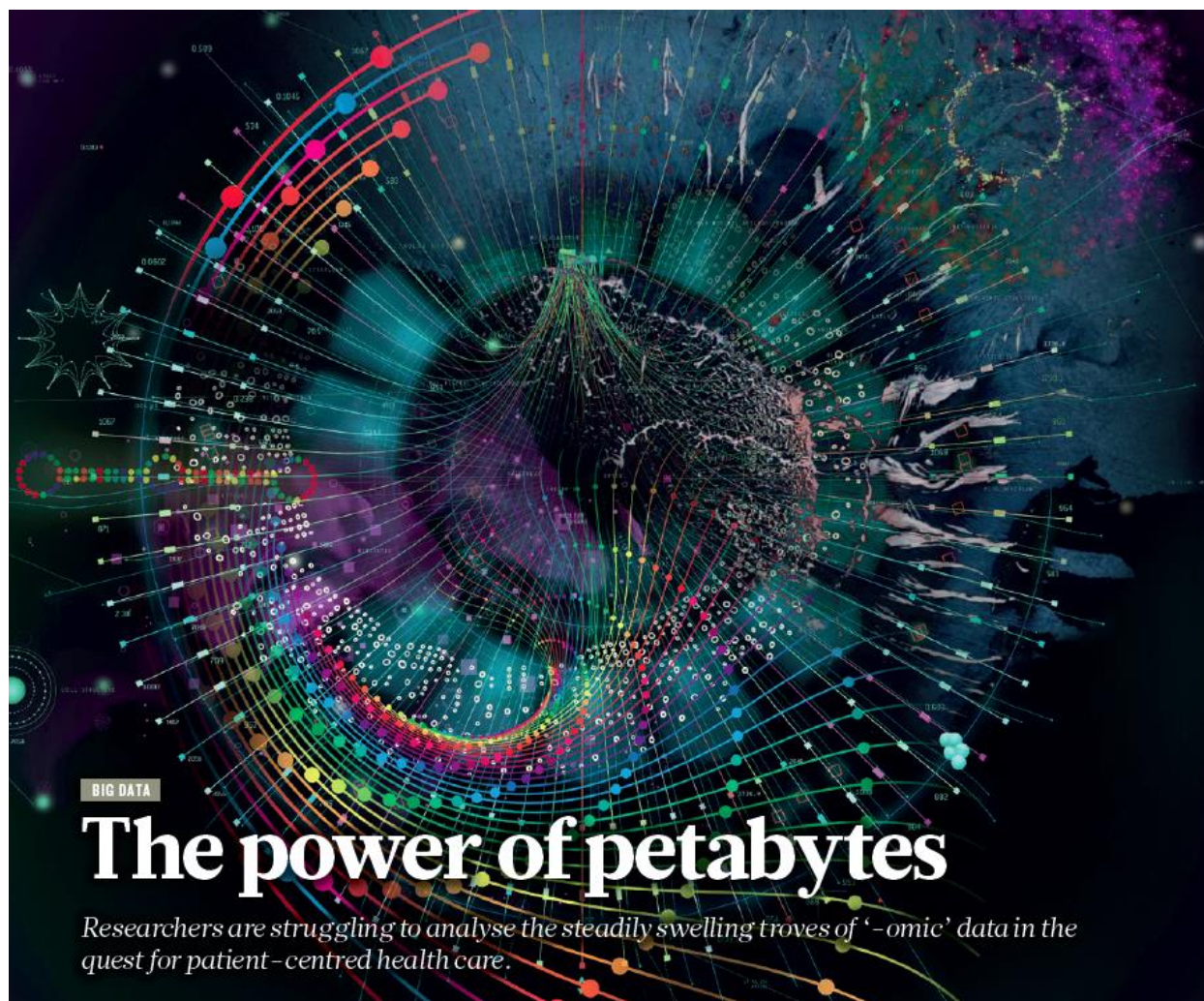
- **NGseq raw imaging data: <span style="color:red">> 2 TB</span> per sample**

- **RNAseq or Exome seq data: <span style="color:red">10 GB</span> per sample (raw data), <span style="color:red">30-50 GB</span> during the processing.**

- **Whole genome seq: <span style="color:red">200 GB</span> per sample (raw data), <span style="color:red">400-600 GB</span> during the processing.**
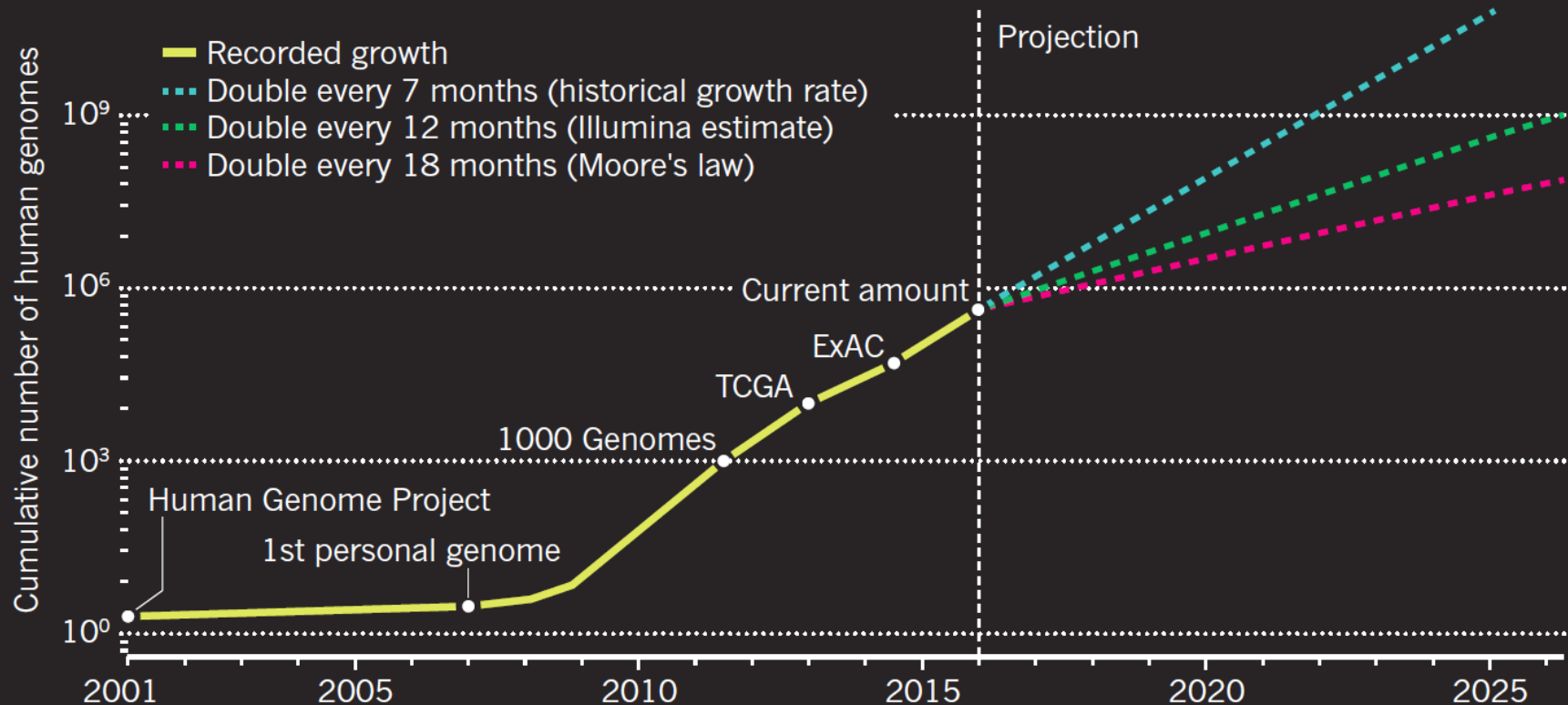
# Data analysis/mining ?

**BIG DATA**

# The power of petabytes

Researchers are struggling to analyse the steadily swelling troves of '-omic' data in the quest for patient-centred health care.

*"You're shooting yourself in the foot if you're collecting data you don't know how to interpret."*

# DNA SEQUENCING SOARS

Human genomes are being sequenced at an ever-increasing rate. The 1000 Genomes Project has aggregated hundreds of genomes; The Cancer Genome Atlas (TGCA) has gathered several thousand; and the Exome Aggregation Consortium (ExAC) has sequenced more than 60,000 exomes. Dotted lines show three possible future growth curves.

Cumulative number of human genomes

- Recorded growth
- Double every 7 months (historical growth rate)
- Double every 12 months (Illumina estimate)
- Double every 18 months (Moore's law)

Projection

$10^9$

$10^6$ — Current amount

ExAC

TCGA

$10^3$

1000 Genomes

Human Genome Project

1st personal genome

$10^0$

2001  2005  2010  2015  2020  2025

- In 2014 the United Kingdom launched the 100,000 Genomes Project, and both the United States (under the Precision Medicine Initiative) and China (in a programme to be run by BGI of Shenzhen) have unveiled plans to analyze genomic data from one million individuals.

- A partnership between Geisinger Health System, based in Danville, Pennsylvania, and biotech firm Regeneron Pharmaceuticals of Tarrytown, New York, for instance, aims to generate sequence data for more than 250,000 people.

# Big data from small data: data-sharing in the 'long tail' of neuroscience

Adam R Ferguson[1], Jessica L Nielson[1], Melissa H Cragin[2], Anita E Bandrowski[3] & Maryann E Martone[3,4]

# Reshaping the cancer clinic

*Big data's war on cancer is still in the early stages, but the front line is advancing.*

The Cancer Genome Atlas, which catalogues cancer mutations, contains some **2.5 million gigabytes of data**. This giant project, run by the US National Institutes of Health, has vastly improved our understanding of various forms of cancer — **but** it holds relatively **little information on the clinical experience** of the patients who supplied the samples.

# BioVU   Vanderbilt DNA Databank

Provides enabling resource for exploration of the relationships among genetic variation, disease susceptibility, and variable drug responses, and represents a key first step in moving the emerging sciences of genomics and pharmacogenomics from research tools to clinical practice. A major goal of the resource is to generate datasets that incorporate de-identified information derived from medical records and genotype information to identify factors that affect disease susceptibility, disease progression, and/or drug response.

**Manager:** Erica Bowton
**E-Mail:** erica.bowton@vanderbilt.edu
**Phone:** (615) 322-1975
**Website:** https://starbrite.vanderbilt.edu/biovu/ (VUnet password required.)

## Vanderbilt BioVU

VANDERBILT **V** UNIVERSITY
MEDICAL CENTER

# What is BioVU?

- The move towards personalized medicine requires very large sample sets for <span style="color:red">discovery</span> and <span style="color:red">validation</span>

- BioVU: biobank intended to support a broad view of biology and enable personalized medicine

- Contains de-identified DNA extracted from <span style="color:red">leftover blood</span> after clinically-indicated testing of Vanderbilt patients who have not opted out

- A major goal of the resource is to generate <span style="color:red">datasets</span> that incorporate de-identified information derived from medical records and genotype information to identify factors that affect <span style="color:red">disease susceptibility</span>, <span style="color:red">disease progression</span>, and/or <span style="color:red">drug response</span>.

# What is the Synthetic Derivative (SD)?

## SCHOOL OF MEDICINE
### VANDERBILT UNIVERSITY

Vanderbilt University

Vanderbilt University Medical Center

VUSM A-Z   VU Directory

## DEPARTMENT OF BIOMEDICAL INFORMATICS

**About Us**  **People**  **Education**  **Research**  **Service**  **FAQs**  **Contacts**

Home
About Us
People
Education
Service
Research
Faculty Development
Alumni
News
Events
FAQs
Contacts

## Synthetic Derivative

The Synthetic Derivative (SD) is the database containing clinical information derived from Vanderbilt's electronic medical record. The SD is a set of records that is no longer linked to the identified medical record from which it is derived and has been altered to the point it no longer closely resembles the original record. The SD can be used as a stand-alone research resource, or can be used in conjunction with BioVU to identify record sets for genome-phenome analysis. The SD interface allows the user to search data extracted from most of the major health information databases at Vanderbilt including StarPanel and the EDW, which is a data warehouse integrating data from EPIC, Medipac and Horizon Export Orders (WIZ). The database contains records for over 2.2 million unique individuals. The search interface allows the user to input basic clinical and demographic information, such as ICD 9 codes, CPT procedure codes, medications, lab values, age and gender and returns de-identified data to the user for review and selection.

DNA samples or genotyping data may be requested after a proposal for the study is received, approved by the BioVU Review Committee and a user agreement is signed. BioVU applications, amendments and data use agreements for BioVU and the Synthetic Derivative are tracked through REDCap databases.

**Faculty Participants:**
Josh Denny
**Aligned Informatics Area:**
Research Informatics

# What is the Synthetic Derivative (SD)?

- **Rich, multi-source <span style="color:red">database</span> of de-identified clinical and demographic data**

- **User Interface tool that can be used for access and analysis**

- **Contains ~<span style="color:red">2.6 million records</span>**

  - **~1 million with detailed longitudinal data**

  - **averaging 100k bytes in size**

  - **an average of 27 codes per record**

# Technology + Policy

**De-identification**

- Derivation of 128-character identifier (RUI) from the MRN generated by Secure Hash Algorithm (SHA-512)

**Date Shift**

- Our algorithm shifts the dates within a record by a time period *(up to 364 days backwards)* that is consistent within each record, but differs *across* records

**Restricted access & continuous oversight**

- IRB approval for study (non-human)
- Data Use Agreement
- Audit logs of all searches and data exports

# Synthetic Derivative (SD)

# Synthetic Derivative Search Features

# VANTAGE

**VANTAGE**

**BioVU**

# VANGARD

## VANGARD

**Van**derbilt Technologies for Advanced **G**enomics **A**nalysis and **R**esearch **D**esign (**VANGARD**) is a new core with administrative oversight from the Office of Research and scientific and technical direction provided by the Vanderbilt Center for Quantitative Sciences. The mission of the core is to consolidate the genomics data pipeline across the university and allow investigators to leverage the opportunities provided by next-generation sequencing and other genomics technologies. VANGARD operates in conjunction with VANTAGE, providing experimental design, quality assessment of data, analysis and results interpretation, and data storage to investigators, while VANTAGE provides technical services with a focus on next-generation sequencing including DNA-seq and RNA-seq. VANGARD also provides biostatistical and bioinformatic support for all genomic experiments that utilize BioVU specimens.

For small-scale projects, VANGARD uses a fee-for-service model which includes basic experimental design and quantitative analysis for genomic data generated by VANTAGE as well as data storage and backup. Large-scale projects and those that require more complex and detailed analysis are handled through a collaborative percent-effort model with VANGARD personnel functioning as research team members.

Dr. Yu Shyr serves as the director of VANGARD, and he maintains close communication with the leadership of VANTAGE to ensure seamless service delivery for genomic research.

Contact
**Genomic design studios.** The first step to a successful study is good experimental design. The VANTAGE/VANGARD team can assist you in designing your genomics experiment to answer the research question of interest. To register for a VANTAGE/VANGARD genomic design studio session, please visit:
http://cqs.mc.vanderbilt.edu/gds

# Vanderbilt ACCRE

# Vanderbilt ACCRE

- **The ACCRE high-performance computing cluster has about 6,000 processor cores.**

- **Compute nodes run 64-bit Linux operating system, with hard drives of 250 GB to 1 TB and dual copper gigabit Ethernet ports.**

- **Each node is monitored via Nagios, with an integrated scheduling system (Moab/Torque) utilized for resource management, scheduling of jobs, and usage tracking.**

- **The home directories of all users are backed up daily to tape.**

# Apple Announced ResearchKit on 4/14/2015

- **ResearchKit™, an open source software framework designed for medical and health research that helps doctors, scientists and other researchers gather data more frequently and more accurately from participants using mobile devices, is now available to researchers and developers.**

- **The first research apps developed using ResearchKit study asthma, breast cancer, cardiovascular disease, diabetes and Parkinson's disease, and have enrolled over 60,000 iPhone users in just the first few weeks of being available on the App Store.**

# iPhone、Apple Watchで脳梗塞を早期発見　国内初の臨床研究、慶大が開始

**iPhoneやApple Watchを活用した臨床研究を慶應義塾大学が開始。不整脈や脳梗塞の早期発見に役立てる考えだ。**

**[ITmedia]**

| 印刷／PDF | | 〈見る | | | | | 💡 通知 |

　慶應義塾大学医学部の研究チームは11月25日、iPhoneやApple Watchのセンサーを活用した臨床研究を国内で初めて開始したと発表した。専用アプリを通じて心拍数や運動能力などを測定し、不整脈・脳梗塞の早期発見につなげるという。iPhoneユーザーであれば、誰でも匿名で参加できる。



専用アプリ「Heart & Brain」を開発

# Apple Announced ResearchKit on 4/14/2015

"Numbers are everything. The more people who contribute their data, the bigger the numbers, the truer the representation of a population, and the more powerful the results. A research platform that allows large amounts of data to be collected and shared — that can only be a positive thing for medical research."

Dr. Eduardo Sanchez, American Heart Association

**Nov 25, 2013**

## 23andMe ordered to halt sales of DNA tests

US regulator seeks information on the safety and effectiveness of the company's analyses.

Sarah Zhang

- After 14 face-to-face and teleconference meetings, hundreds of email exchanges, and dozens of written communications, FDA orders 23andMe to halt sales of DNA test kit.

- 23andMe Uses a sample of a user's saliva

- Claims to identify up to 254 diseases and other medical conditions

- 23andMe has submitted to the FDA twice for review as a medical device but has failed to satisfy the FDA's concerns

# 23andMe is back in the genetic testing business with FDA approval Oct, 2015

- **In 2013: Assessments on <span style="color:red">254</span> diseases and conditions for $99**

- **In 2015: Assessments on <span style="color:red">36</span> inherited disorders for $199**

- **More than one million subscribers in the company's database**

- **The company also just announced a <span style="color:red">new drug discovery and development venture</span>.**

# Cancer Precision Medicine

Spotlight

# Data Scientist:
## *The Sexiest Job of the 21st Century*

A new role is fast gaining prominence in organizations: that of the data scientist. Data scientists are the people who understand how to fish out answers to important business questions from today's tsunami of unstructured information. As companies rush to capitalize on the potential of big data, the largest constraint many face is the scarcity of this special talent.

The shortage of data scientists is becoming a serious constraint in some sectors.

Data scientists today are akin to the Wall Street "quants" of the 1980s and 1990s.

Ben Chams - Fotolia

# NGS – Data Analysis

**Table 4 Computational tools for cancer transcriptomics**

| Category | Program | URL | ref |
| --- | --- | --- | --- |
| Spliced alignment | TopHat | http://tophat.cbcb.umd.edu/ | [61,69] |
| | MapSplice | http://www.netlab.uky.edu/p/bioinfo/MapSplice | [62] |
| | SpliceMap | http://www.stanford.edu/group/wonglab/SpliceMap/ | [63] |
| | GSNAP | http://research-pub.gene.com/gmap/ | [64] |
| | STAR | http://gingeraslab.cshl.edu/STAR/ | [65] |
| Differential expression | CuffDiff | http://cufflinks.cbcb.umd.edu/ | [68,69] |
| | EdgeR | http://www.bioconductor.org/packages/2.11/bioc/html/edgeR.html | [67] |
| | DESeq | http://www-huber.embl.de/users/anders/DESeq/ | [66] |
| | Myrna | http://bowtie-bio.sourceforge.net/myrna/index.shtml | [81] |
| Alternative splicing | CuffDiff | http://cufflinks.cbcb.umd.edu/ | [68,69] |
| | MISO | http://genes.mit.edu/burgelab/miso/ | [71] |
| | DEXseq | http://watson.nci.nih.gov/bioc_mirror/packages/2.9/bioc/html/DEXSeq.html | [82] |
| | Alexa-seq | http://www.alexaplatform.org/alexa_seq/ | [70] |
| Gene fusion | SOAPfusion | http://soap.genomics.org.cn/SOAPfusion.html | |
| | TopHat-Fusion | http://tophat.cbcb.umd.edu/fusion_index.html | [72] |
| | BreakFusion | http://bioinformatics.mdanderson.org/main/BreakFusion | [73] |
| | FusionHunter | http://bioen-compbio.bioen.illinois.edu/FusionHunter/ | [74] |
| | deFuse | http://sourceforge.net/apps/mediawiki/defuse/ | [75] |
| | FusionAnalyser | http://www.ilte-cml.org/FusionAnalyser/ | [76] |

**BMC Bioinformatics**

# Sample size calculation based on exact test for assessing differential expression analysis in RNA-seq data

Chung-I Li[1,3], Pei-Fang Su[2,3] and Yu Shyr[3*]

## NIH Public Access
### Author Manuscript

### Sample size calculation for differential expression analysis of RNA-seq data under Poisson distribution

**Chung-I Li**, **Pei-Fang Su**, **Yan Guo**, and **Yu Shyr**
Center for Quantitative Sciences, Vanderbilt University, 571 Preston Building Nashville, TN, USA

## Biometrics & Biostatistics

# Sample Size Calculation of RNA-sequencing Experiment-A Simulation-Based Approach of TCGA Data

# Biometrics & Biostatistics

**Research Article**                                                                                                           **Open Access**

# Sample Size Calculation of RNA-sequencing Experiment-A Simulation-Based Approach of TCGA Data
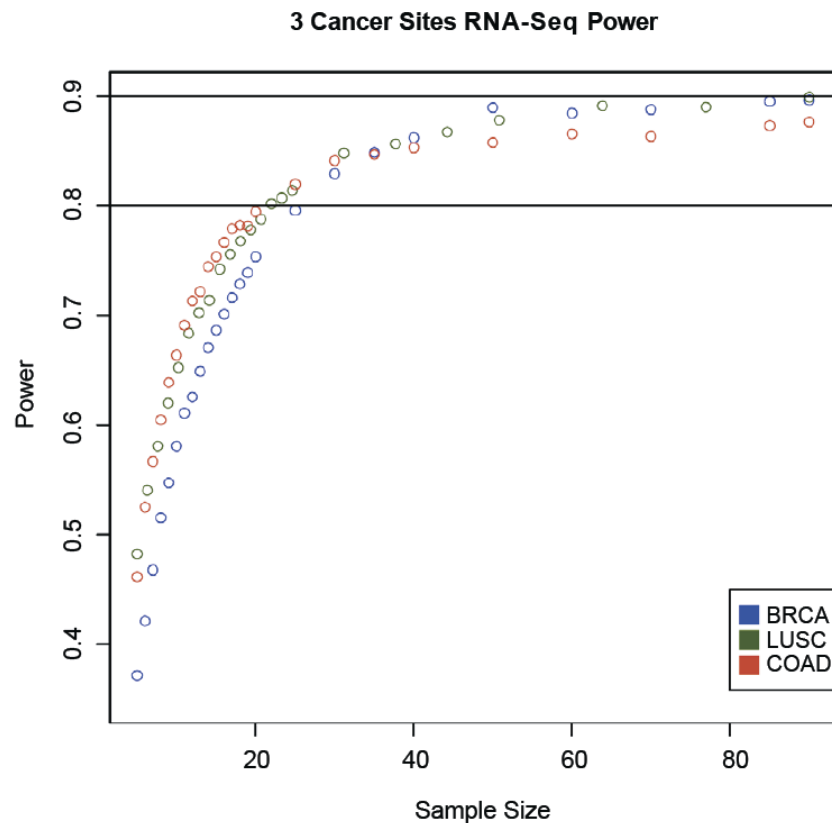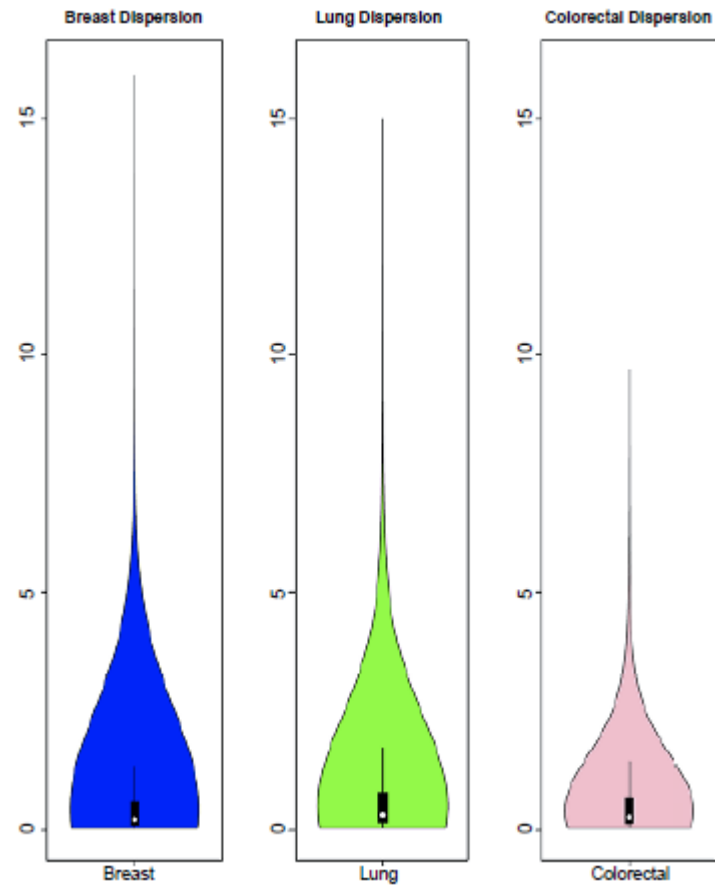


**Figure 3:** Scatterplot of the samples size and power (minimum reads=5 and FDR=0.05).

# Illumina human exome genotyping array clustering and quality control

Yan Guo[1], Jing He[2], Shilin Zhao[1], Hui Wu[1], Xue Zhong[1], Quanhu Sheng[1], David C Samuels[3], Yu Shyr[1] & Jirong Long[2]

[1]Center for Quantitative Sciences, Vanderbilt University, Nashville Tennessee, USA. [2]Vanderbilt Epidemiology Center, Vanderbilt University, Nashville Tennessee, USA. [3]Center for Human Genetics Research, Vanderbilt University, Nashville, Tennessee, USA. Correspondence should be addressed to Y.G. (yan.guo@vanderbilt.edu).
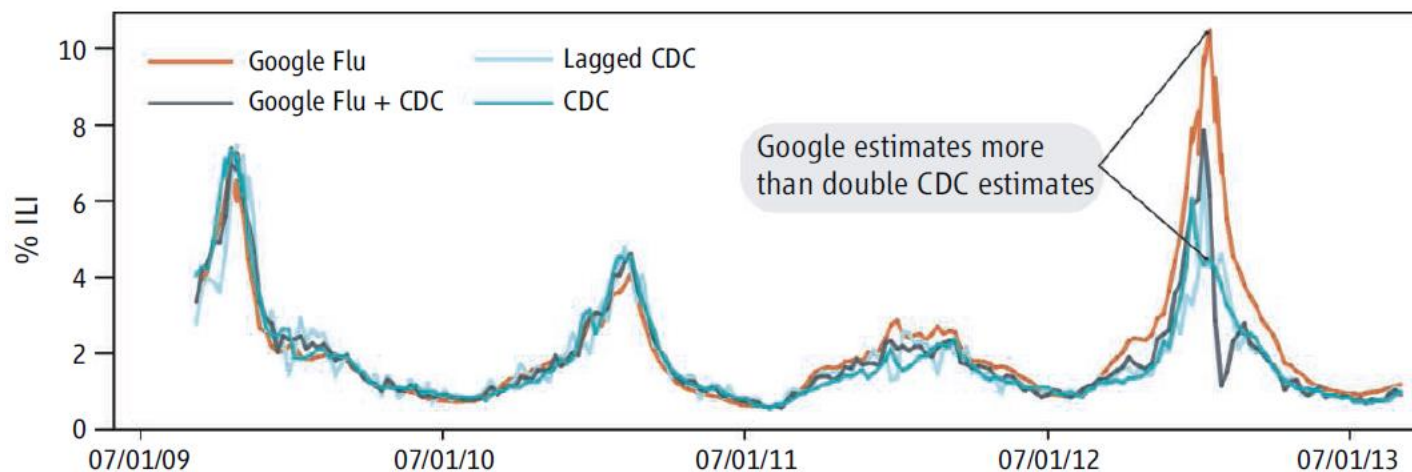
**GenomeStudio section**
- Loading the data into GenomeStudio
- Performing automatic clustering
- QC on SNPs located in a haploid genome
- QC based on GenTrain score
- QC based on cluster separation
- QC based on Mendelian error and replication error
- QC based on other criteria
- Calling rare SNPs
- Final filtering
- Exporting from GenomeStudio

**Post-GenomeStudio section**
- Converting all SNPs to the forward strand
- Checking for gender mismatch
- Checking for race mismatch
- Checking for relatedness
- Checking for Hardy-Weinberg equilibrium (HWE) outliers
- Checking for heterozygosity outliers
- Checking consistency between exome chip genotype and 1000 Genomes Project[17] or HapMap[18] genotype
- Checking for minor allele frequency (MAF) consistency between exome chip and 1000 Genomes Project genotypes
- Checking for batch effects

# The Parable of Google Flu: Traps in Big Data Analysis

David Lazer,[1,2]* Ryan Kennedy,[1,3,4] Gary King,[3] Alessandro Vespignani[5,6,3]



- This should have been a warning that the big data were **over-fitting** the small number of cases—a standard concern in data analysis.
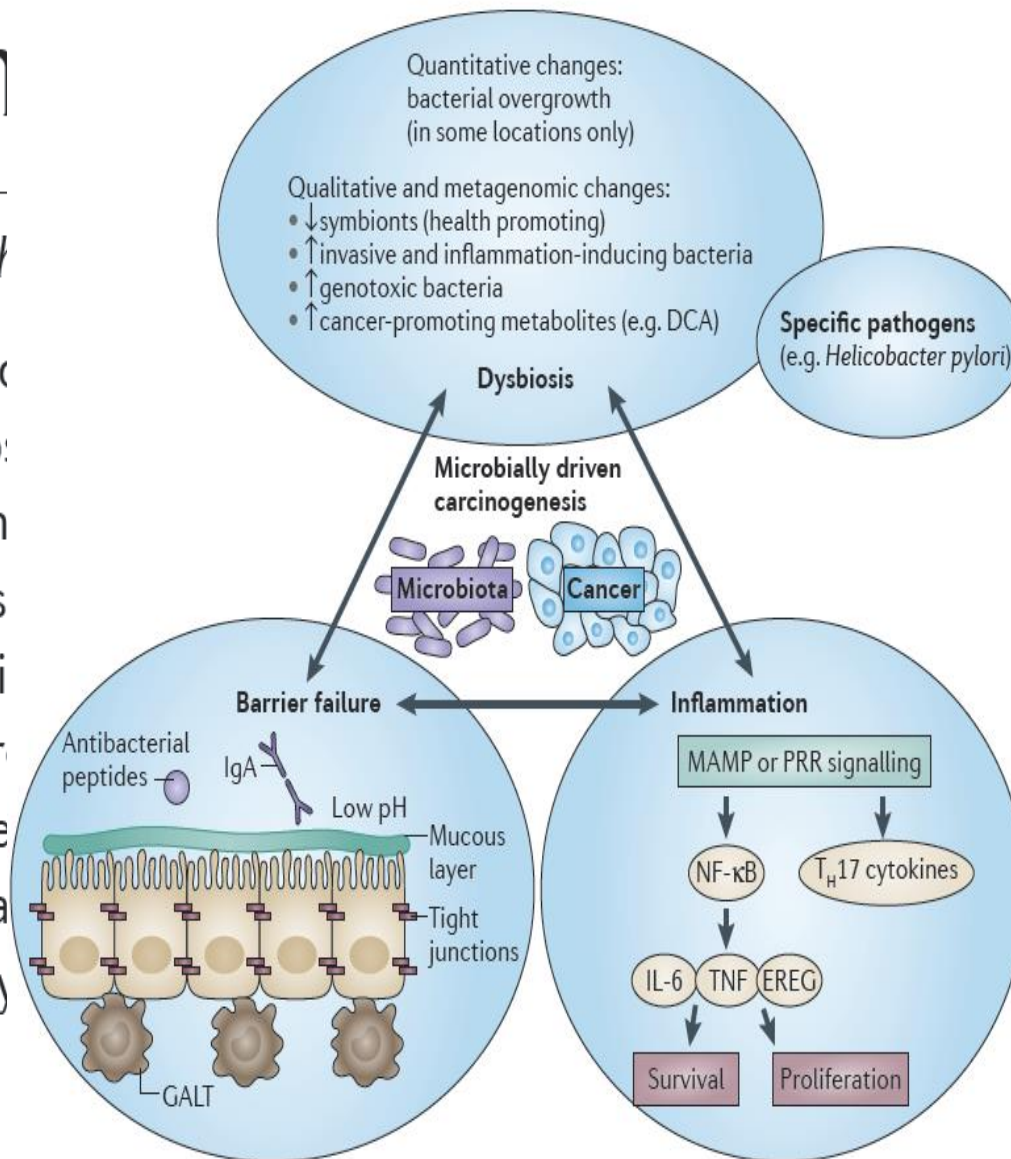
**Microbiome and PheWAS**

# The m

*Robert F. Sc*

Abstract | Mic ... ch symbiotic relationships ... fe. However, defects in th ... sensing and homeostasis ... ntal changes (infection, di ... nd promote disease. Incr ... icrobiota in carcinogene ... he bacterial microbiota a ... es, dysbiosis, genotoxicity ... e for cancer prevention.

Quantitative changes: bacterial overgrowth (in some locations only)

Qualitative and metagenomic changes:
- ↓symbionts (health promoting)
- ↑invasive and inflammation-inducing bacteria
- ↑genotoxic bacteria
- ↑cancer-promoting metabolites (e.g. DCA)

**Dysbiosis**

**Specific pathogens** (e.g. *Helicobacter pylori*)

**Microbially driven carcinogenesis**

Microbiota | Cancer

**Barrier failure**

Antibacterial peptides — IgA — Low pH — Mucous layer — Tight junctions — GALT

**Inflammation**

MAMP or PRR signalling

NF-κB | T$_H$17 cytokines

IL-6 | TNF | EREG

Survival | Proliferation

# Emerging roles of the microbiome

Research on the microbiome is an emerging science. Recent work suggests the benefits derived by humans from their microbiotas may have profound consequences for health.

Roles include
- Food digestion and nutrition
- Regulation of metabolism
- Processing and detoxifying environmental chemicals
- Development and regulation of the immune system
- Prevention of invasion and growth of pathogens
- Role in carcinogenesis: cancer susceptibility and progression
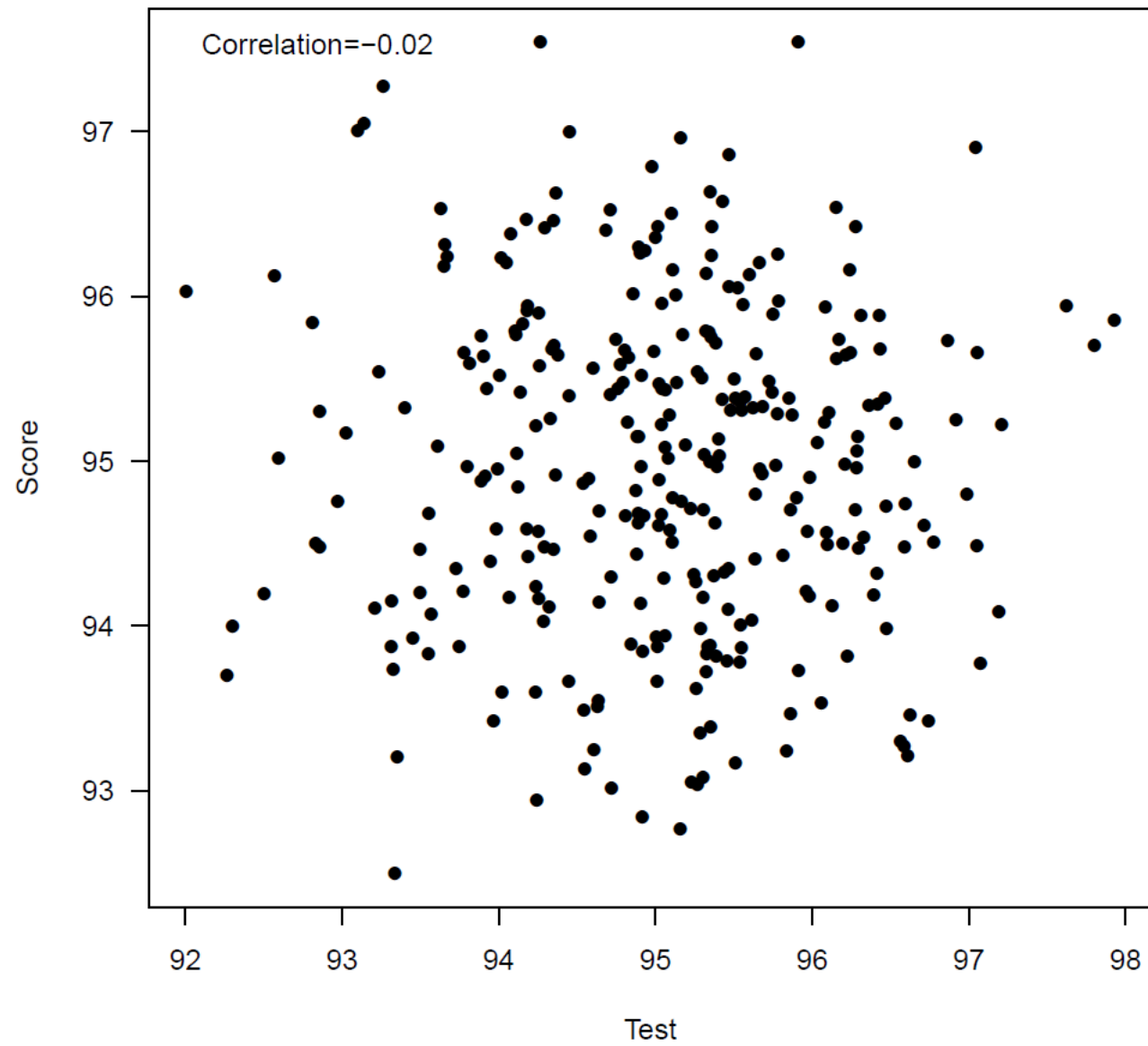- Biological control of certain diseases and deficiencies
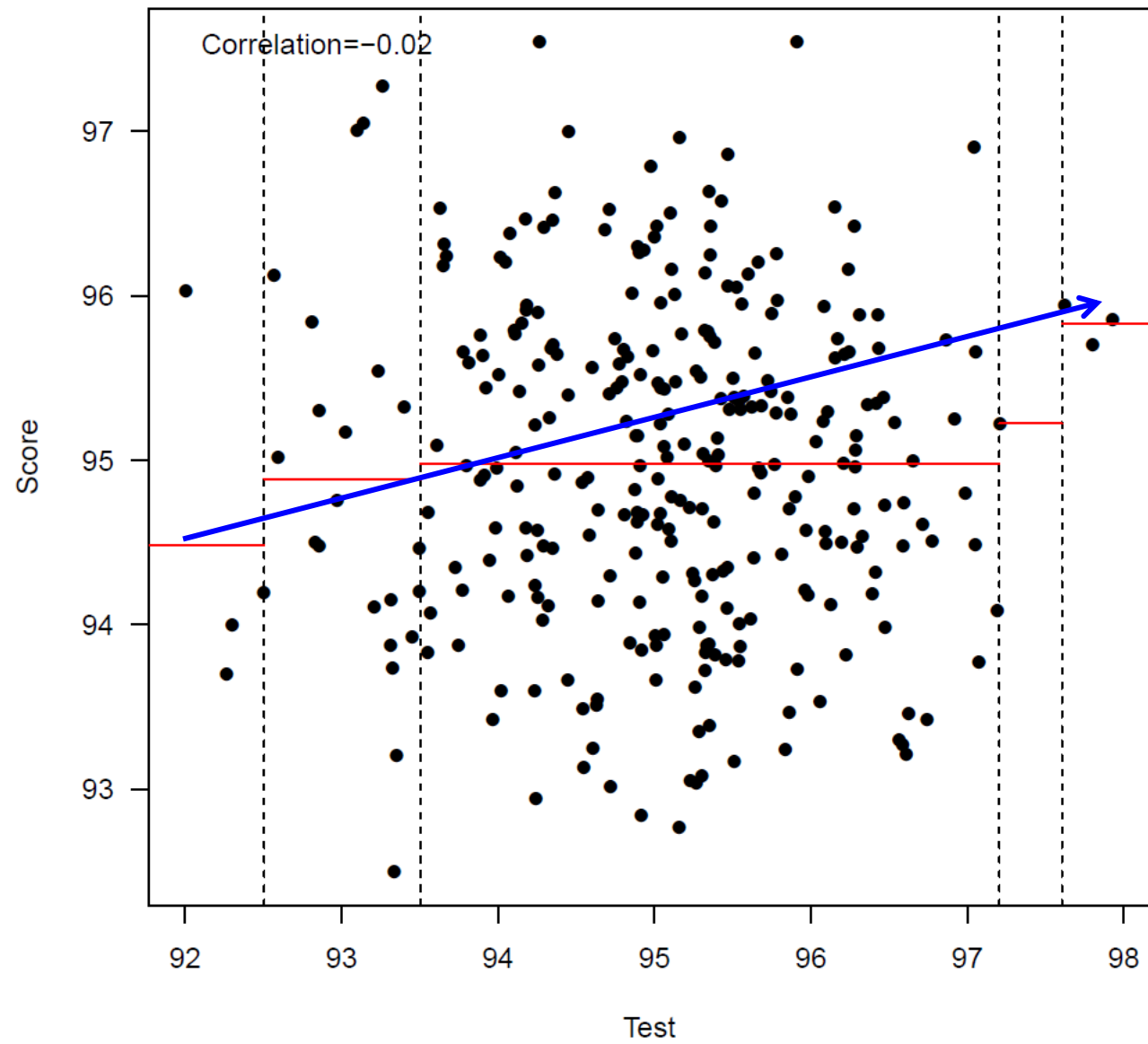
# PheWAS

- **PheWAS explores the association between a number of common genetic variations and a wide variety of phenotypes**

- **The combination of the extensive collection of de-identified medical records in the Synthetic Derivative and the genomic information in BioVU is ideal for PheWAS**
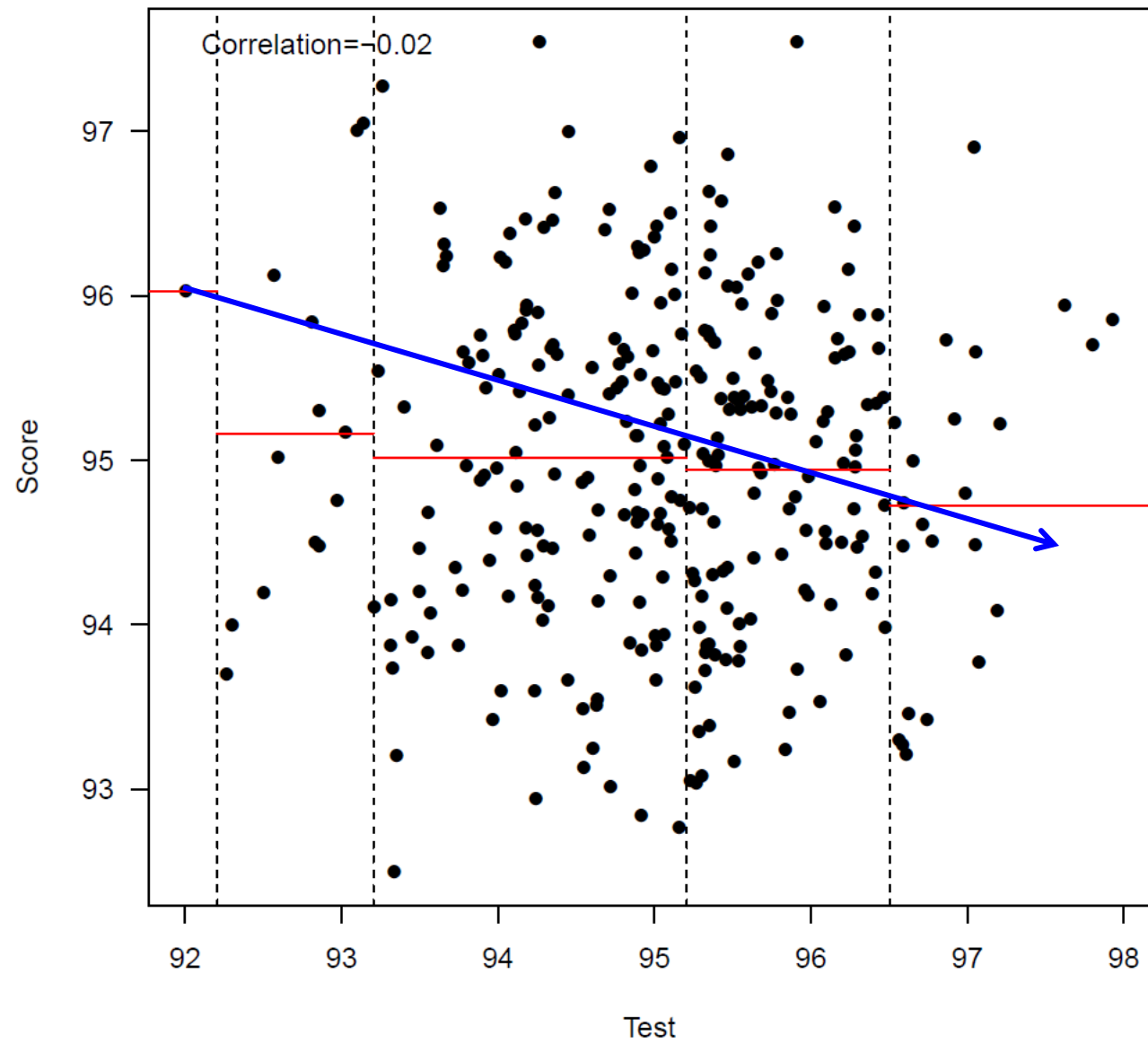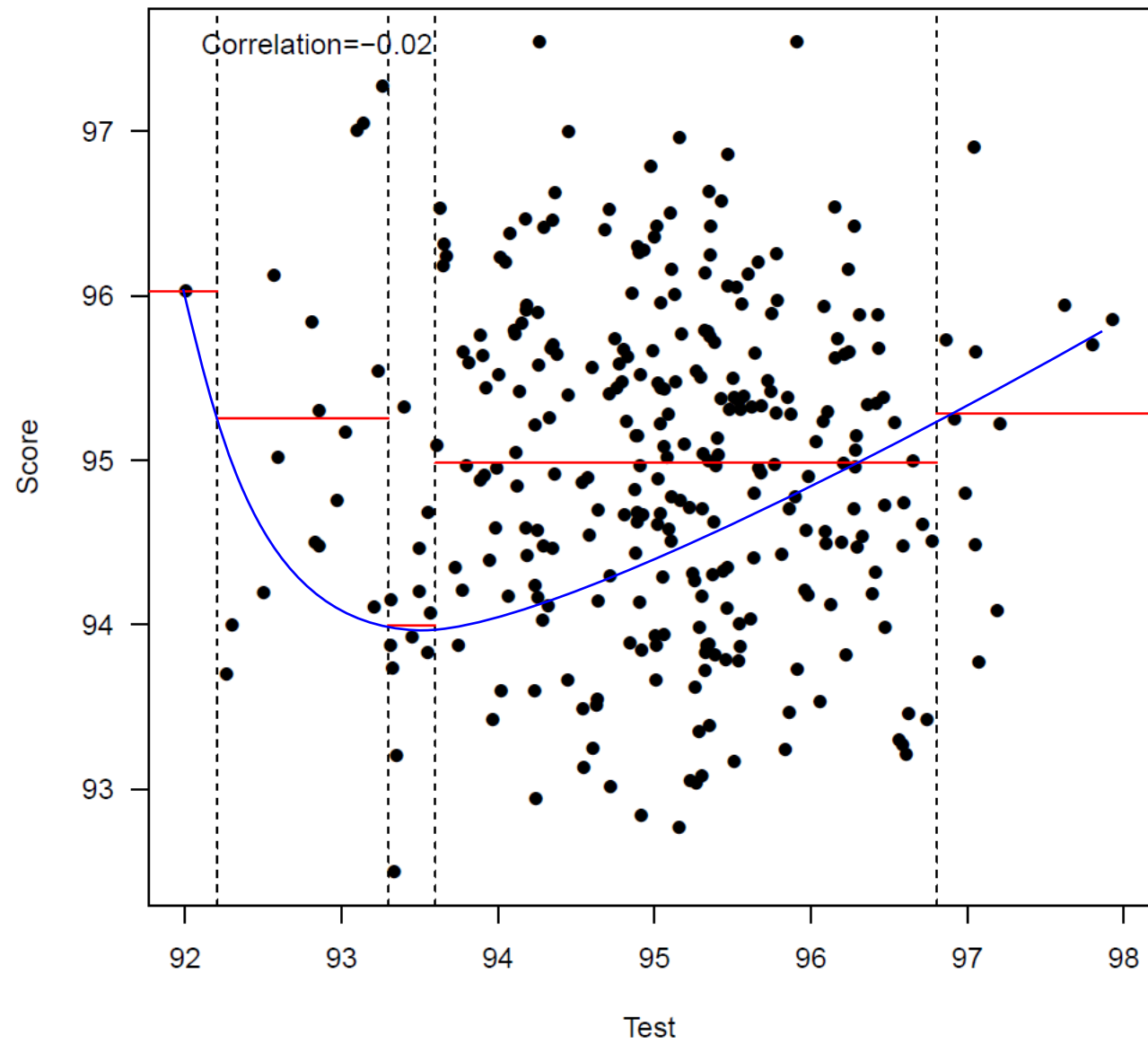
# PheWAS

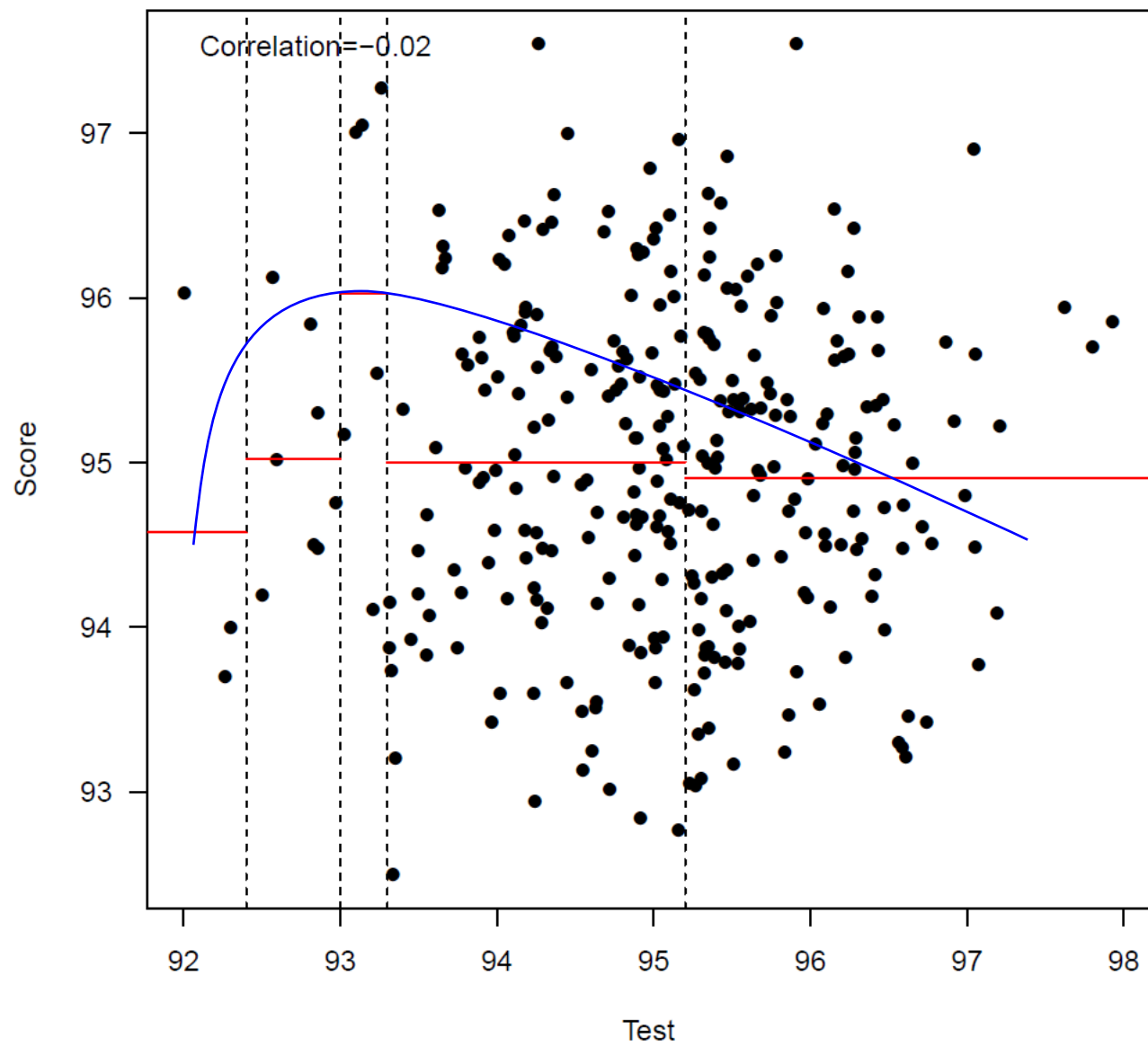**PheWAS: demonstrating the feasibility of a phenome-wide scan to discover gene–disease associations**

# END

# Questions