# "Quantifying Your Superorganism Body Using Big Data Supercomputing"

**ACM International Workshop on Big Data in Life Sciences**

**BigLS 2014**

**Newport Beach, CA**

**September 20, 2014**

**Dr. Larry Smarr**

**Director, California Institute for Telecommunications and Information Technology**

**Harry E. Gruber Professor,**

**Dept. of Computer Science and Engineering**

**Jacobs School of Engineering, UCSD**

**http://lsmarr.calit2.net**

1

# Abstract

The human body is host to 100 trillion microorganisms, ten times the number of cells in the human body and these microbes contain 100 times the number of DNA genes that our human DNA does. The microbial component of this "superorganism" is comprised of hundreds of species spread over many taxonomic phyla. The human immune system is tightly coupled with this microbial ecology and in cases of autoimmune disease, both the immune system and the microbial ecology can have dynamic excursions far from normal. Our research starts with trillions of DNA bases, produced by Illumina Next Generation sequencers, of the human gut microbial DNA taken from my own body, as well as from hundreds of people sequenced under the NIH Human Microbiome Project. To decode the details of the microbial ecology we feed this data into parallel supercomputers, running sophisticated bioinformatics software pipelines. We then use Calit2/SDSC designed Big Data PCs to manage the data and drive innovative scalable visualization systems to examine the complexities of the changing human gut microbial ecology in health and disease. Finally, I will show how advanced data analytics tools find patterns in the resulting microbial distribution data that suggest new hypotheses for clinical application.

**Newsweek**

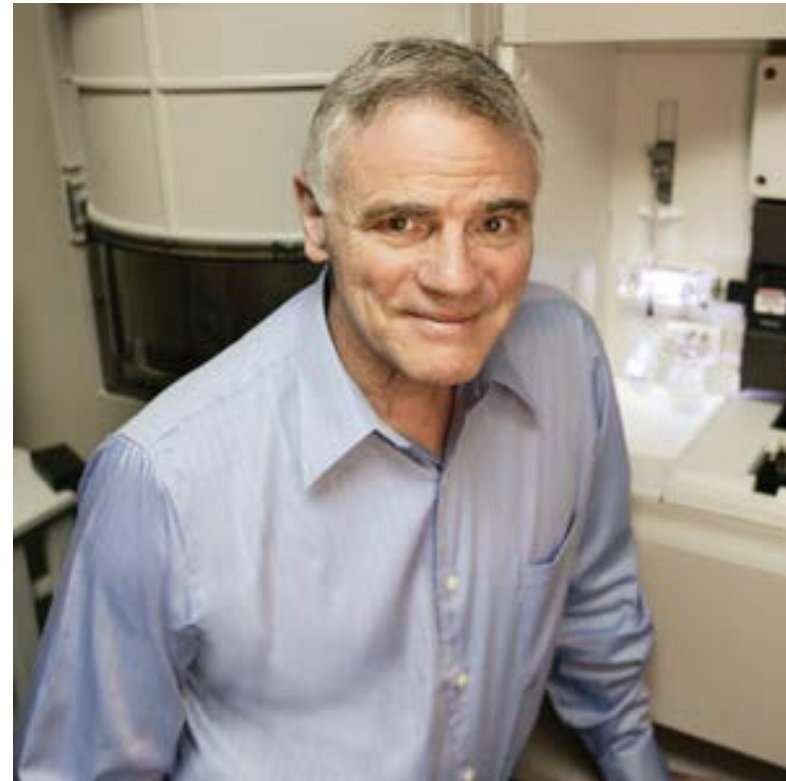# A Doctor's Vision of the Future of Medicine

Leroy Hood
NEWSWEEK
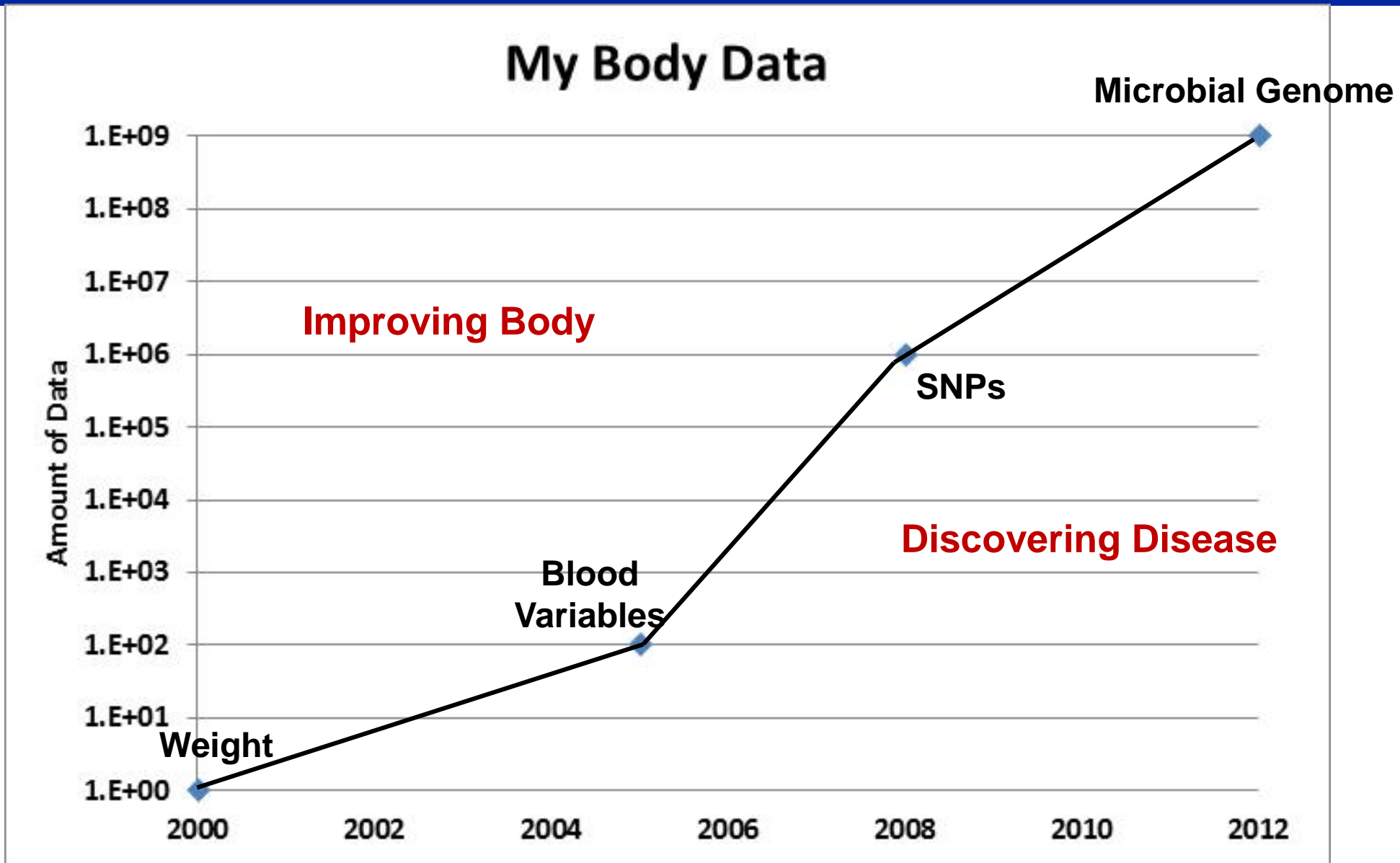
From the magazine issue dated Jul 13, 2009

**PIONEER 100**

HUNDRED PERSON WELLNESS PROJECT
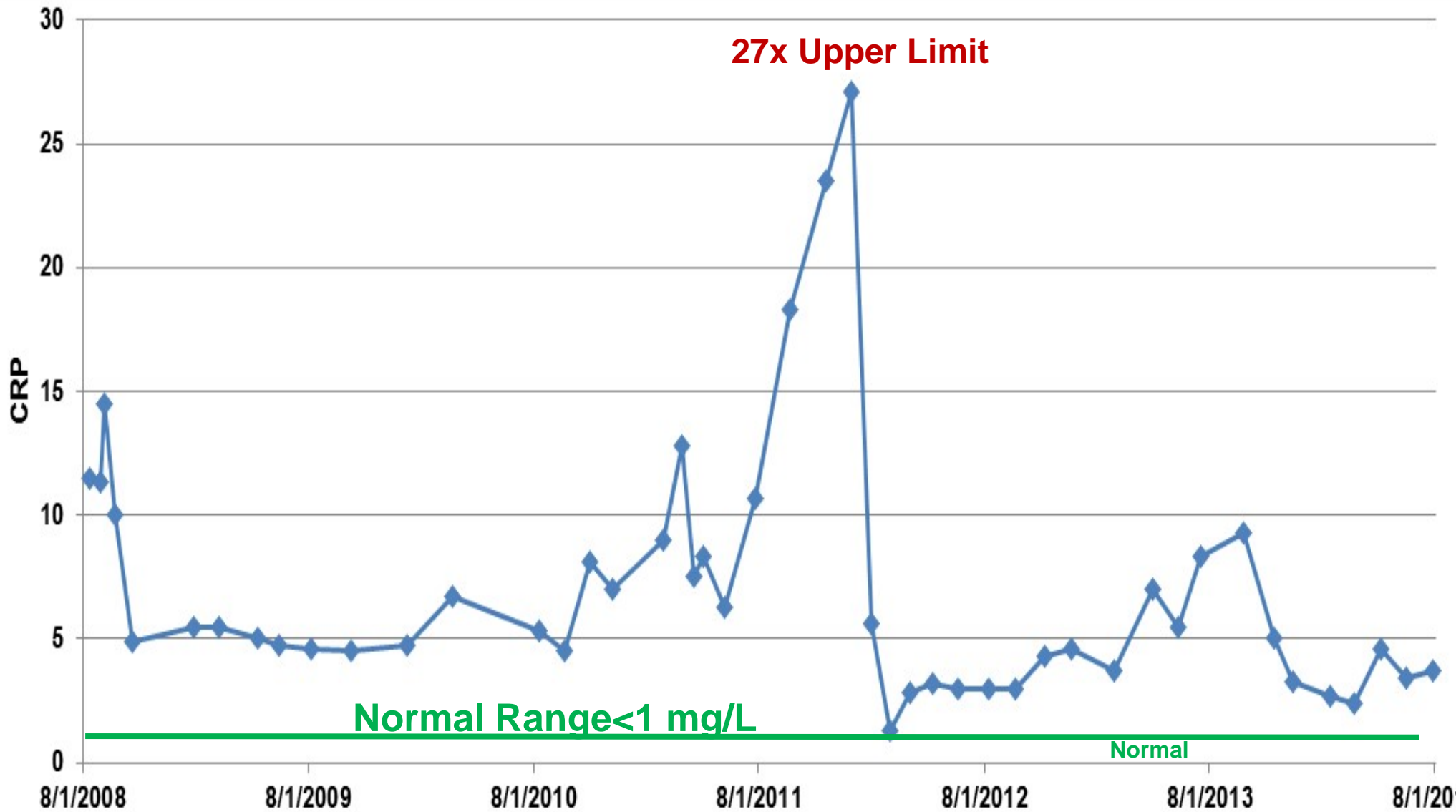
Institute for Systems Biology

**Will Grow to 1000, then 10,000**

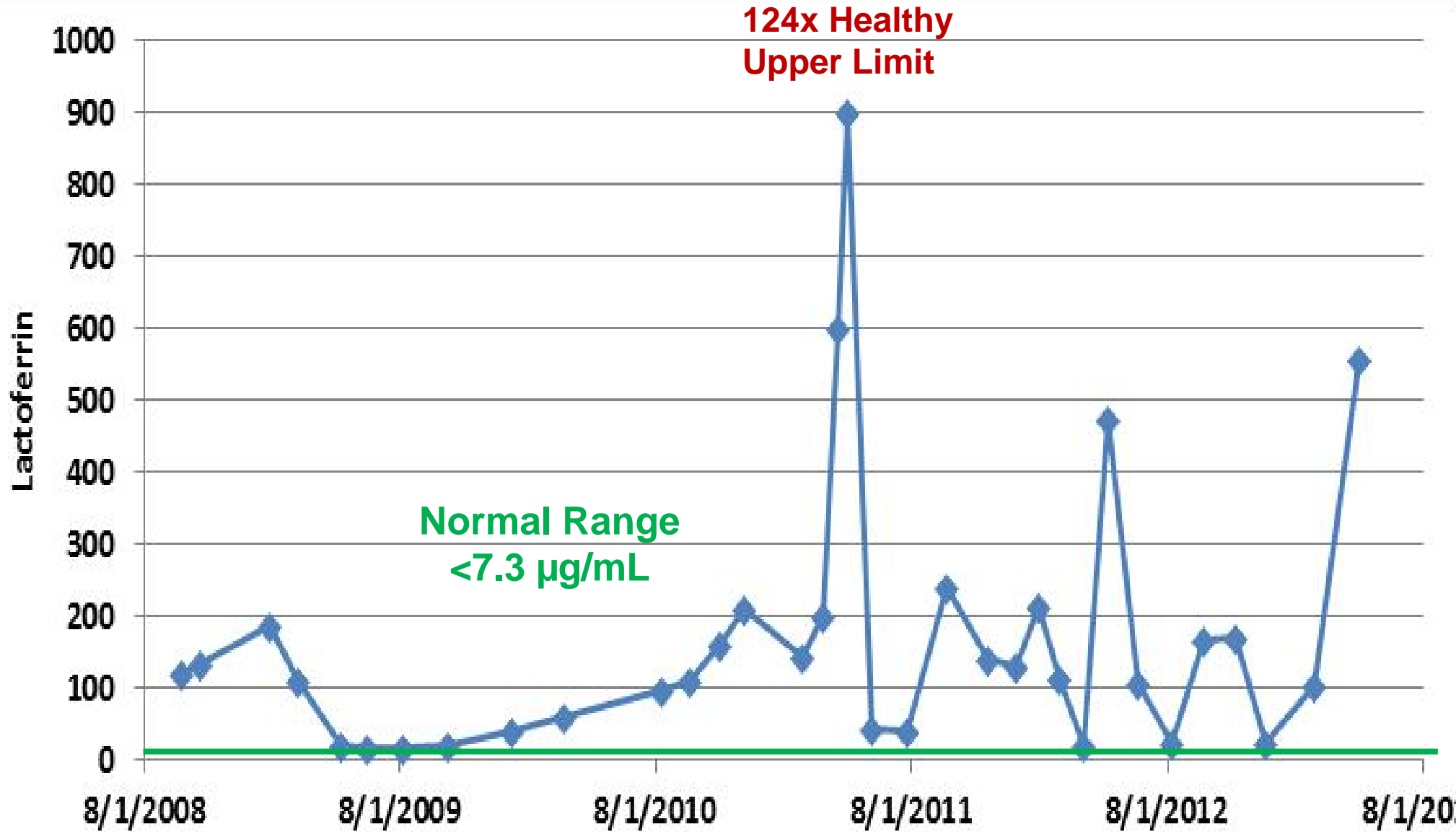# From One to a Billion Data Points Defining Me: The Exponential Rise in Body Data in Just One Decade

# Visualizing Time Series of 150 LS Blood and Stool Variables, Each Over 5-10 Years

**Calit2 64 megapixel VROOM**

# One of My Blood Measurements
# Was Far Out of Range--Indicating Chronic Inflammation



**27x Upper Limit**

**Normal Range<1 mg/L**

Normal

CRP

**Complex Reactive Protein (CRP) is a Blood Biomarker for Detecting Presence of Inflammation**

**Stool Samples Revealed Episodic Autoimmune Response**

124x Healthy Upper Limit

Normal Range <7.3 µg/mL

Lactoferrin is an Antibacteria Glycoprotein Shed from Attacking WBC Neutrophils

# High Lactoferrin Biomarker Led Me to Hypothesis I Had Inflammatory Bowel Disease (IBD)

**IBD is an Autoimmune Disease Which Comes in Two Subtypes: Crohn's and Ulcerative Colitis**



Scand J Gastroenterol. 42, 1440-4 (2007)

My Values May 2011 →

My Values 2009-10 →

**High Level of Calprotectin Confirmed Hypothesis**

# Why Did I Have an Autoimmune Disease like IBD?

Despite **decades of research**,
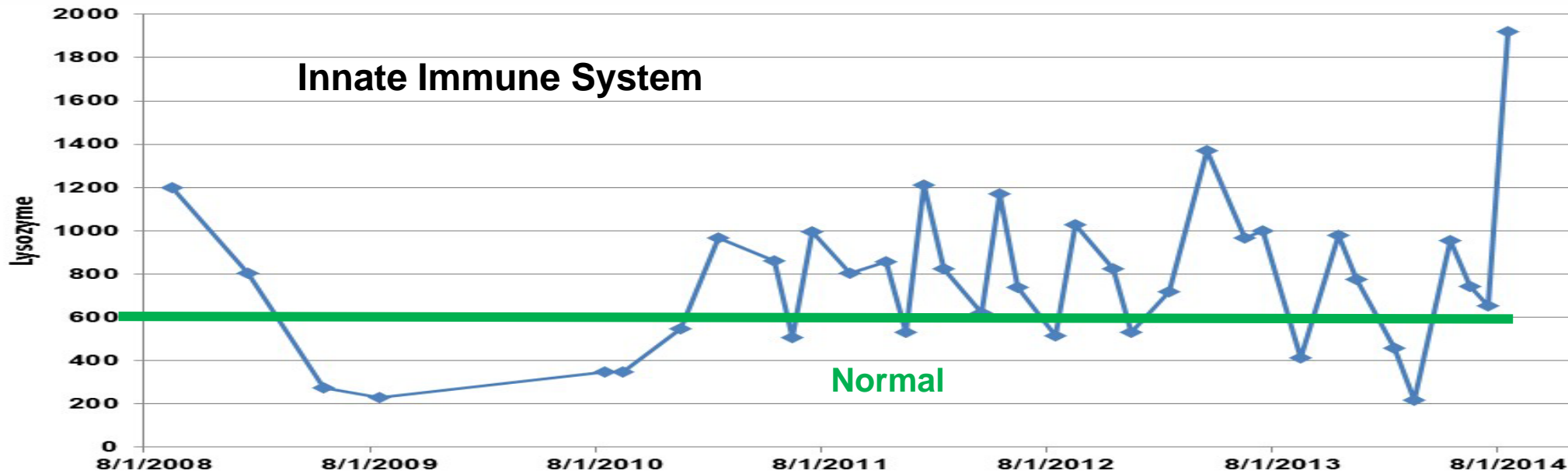the etiology of Crohn's disease
**remains unknown**.
Its pathogenesis may involve
a **complex interplay** between
**host genetics**,
**immune dysfunction**,
and **microbial** or environmental factors.
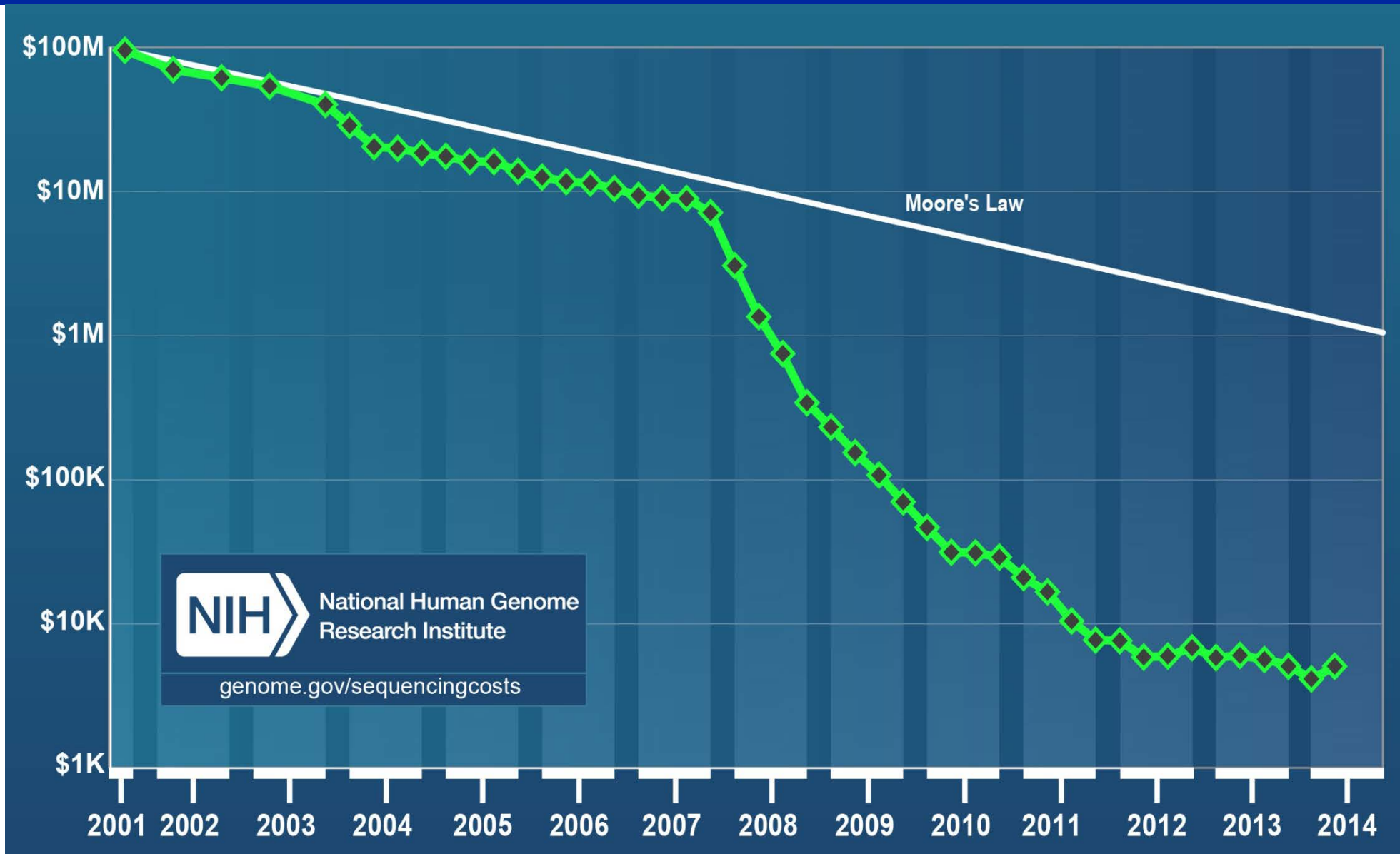*--The Role of Microbes in Crohn's Disease*

## So I Set Out to Quantify All Three!

# Fine Time-Resolution Sampling Reveals Dynamical Innate and Adaptive Immune Dysfunction

# I Found I Had One of the Earliest Known SNPs Associated with Crohn's Disease



From www.23andme.com

**Interleukin-23 Receptor Gene — 80% Higher Risk of Pro-inflammatory Immune Response**

SNPs Associated with CD

I am an Advisor to 23andme
Who Are Seeking
10,000 Volunteers with IBD
to Determine SNP Distribution
to Stratify Disease Spectrum

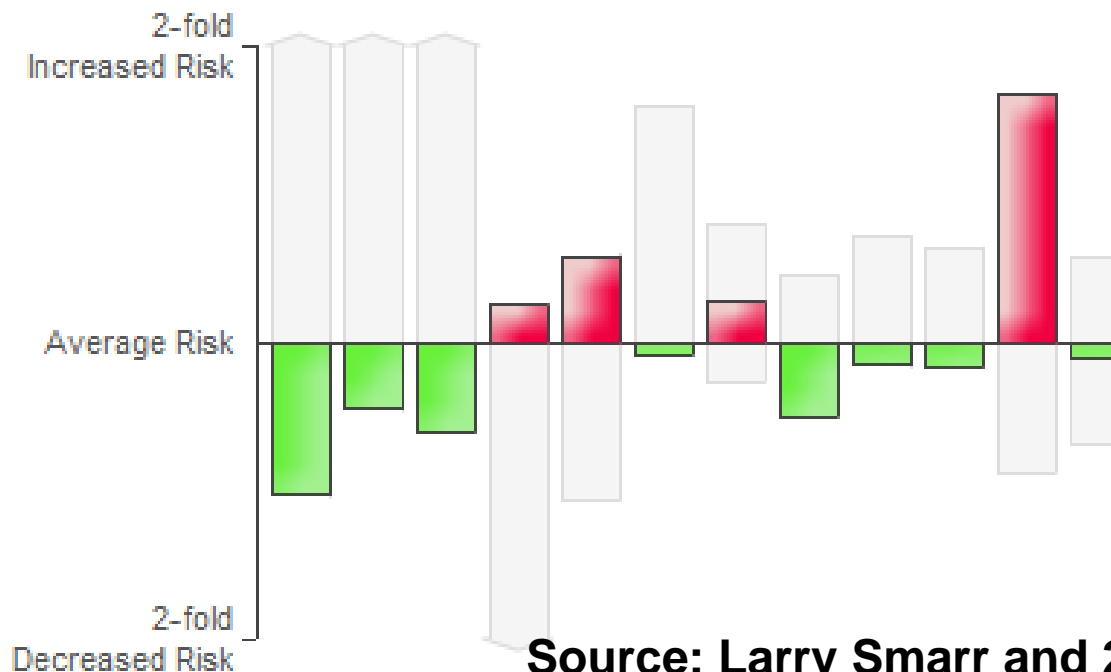News and Views

IL-23: a master regulator in Crohn disease

Markus F Neurath[1]

# There Is Likely a Correlation Between CD SNPs and Where and When the Disease Manifests



**Subject with Ileal Crohn's (ICD)**

**NOD2 (1) rs2066844**

**Female CD Onset At 20-Years Old**

**Subject with Colon Crohn's (CCD)**
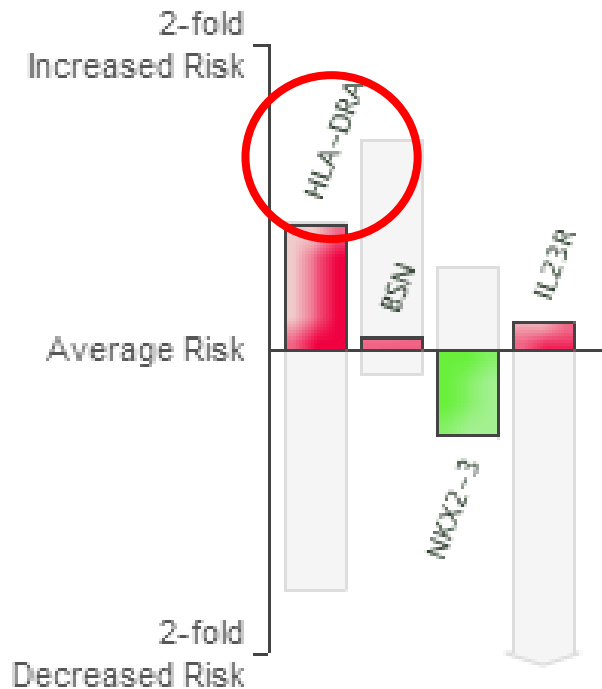
**Il-23R rs1004819**

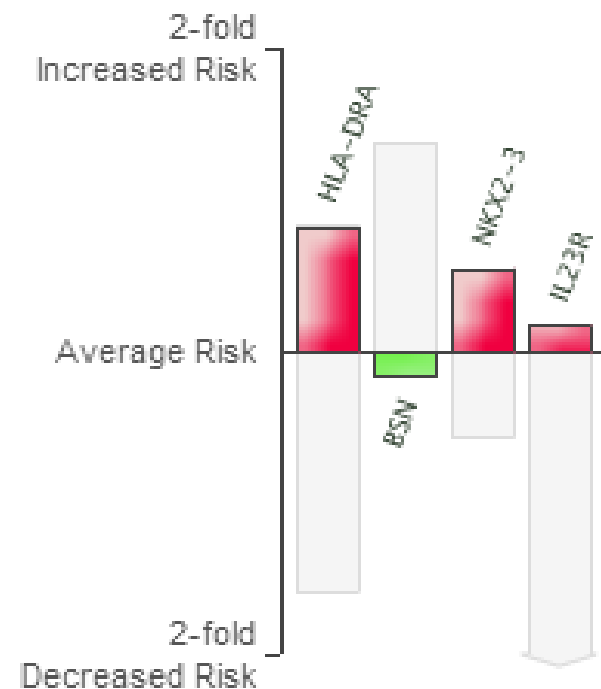**Me-Male CD Onset At 60-Years Old**

**Source: Larry Smarr and 23andme**

# I Also Had an Increased Risk for Ulcerative Colitis, But a SNP that is Also Associated with Colonic CD



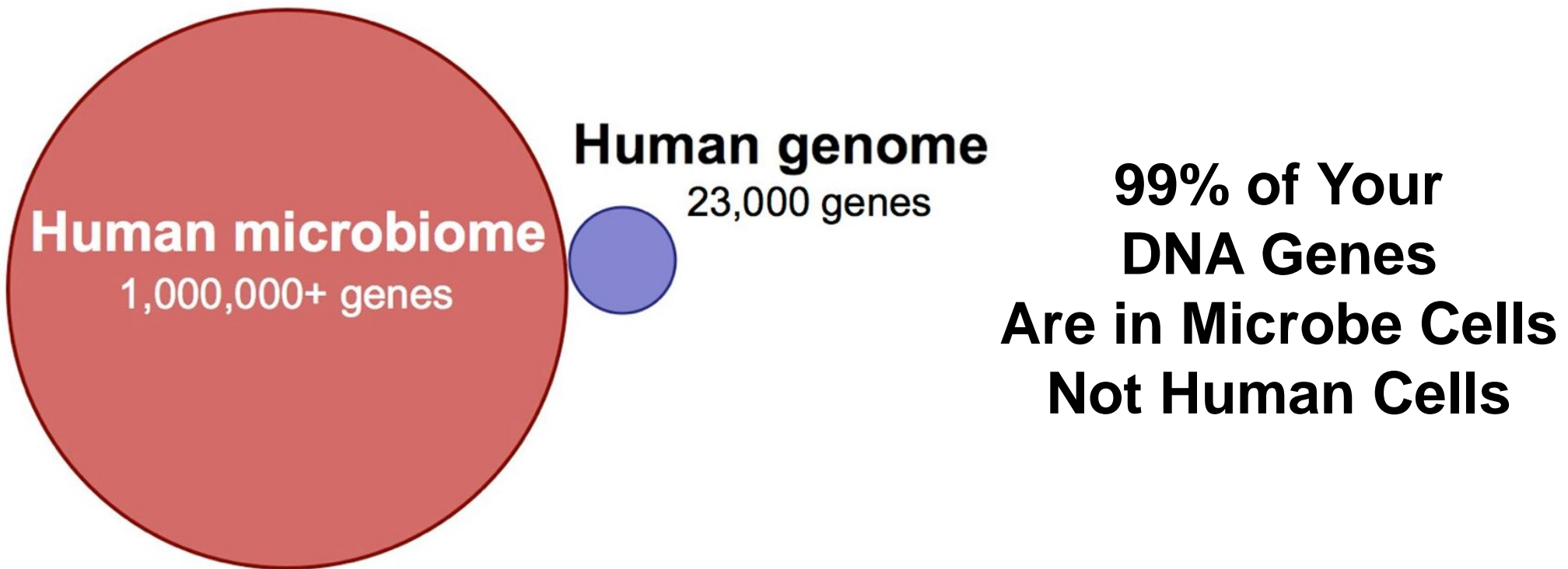**I Have a 33% Increased Risk for Ulcerative Colitis HLA-DRA (rs2395185)**

**I Have the Same Level of HLA-DRA Increased Risk as Another Male Who Has Had Ulcerative Colitis for 20 Years**

"Our results suggest that at least for the SNPs investigated [including HLA-DRA], colonic CD and UC have common genetic basis."
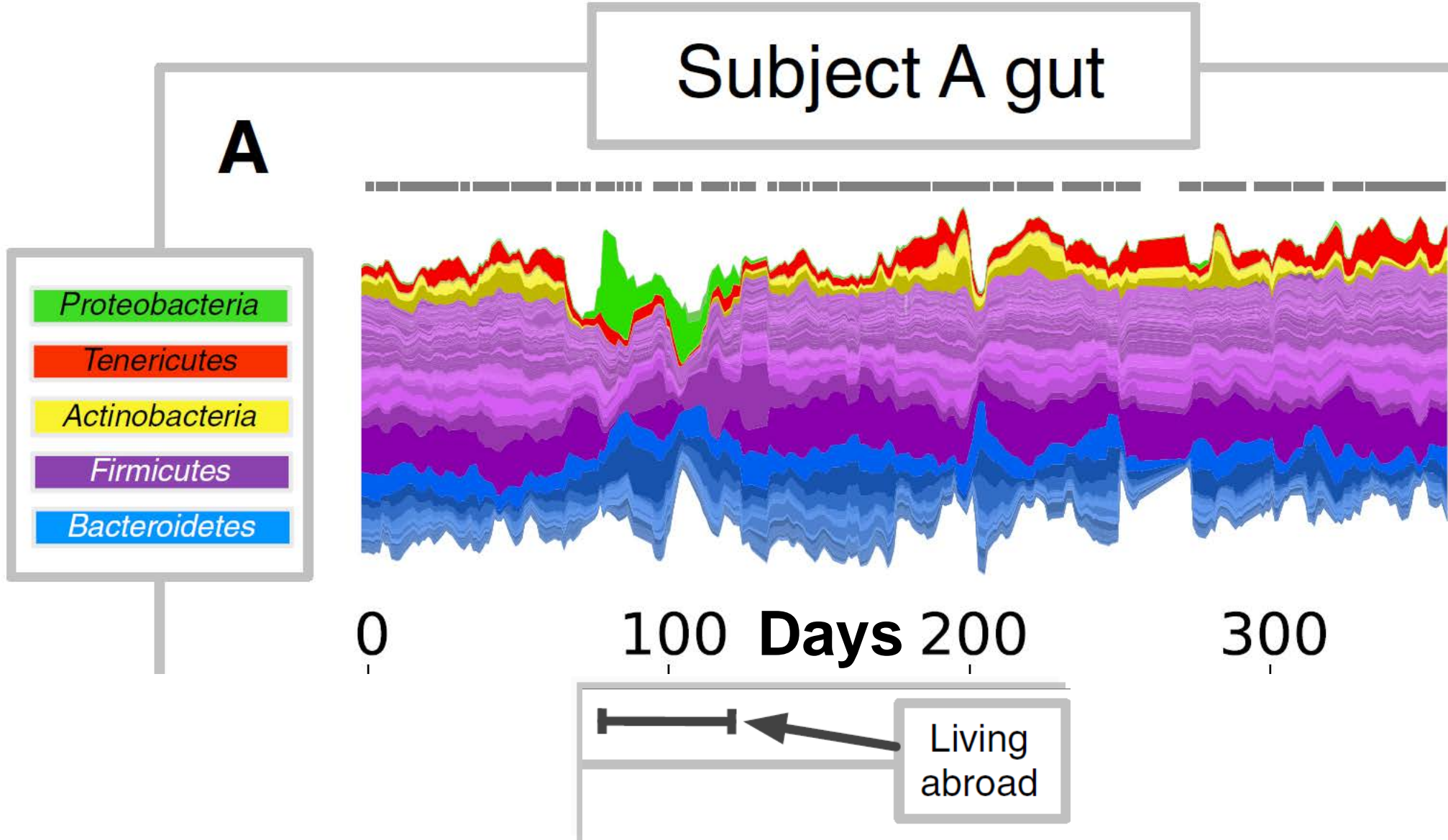-Waterman, et al., IBD 17, 1936-42 (2011)

# Now I am Observing the 100 Trillion Non-Human Cells in My Body

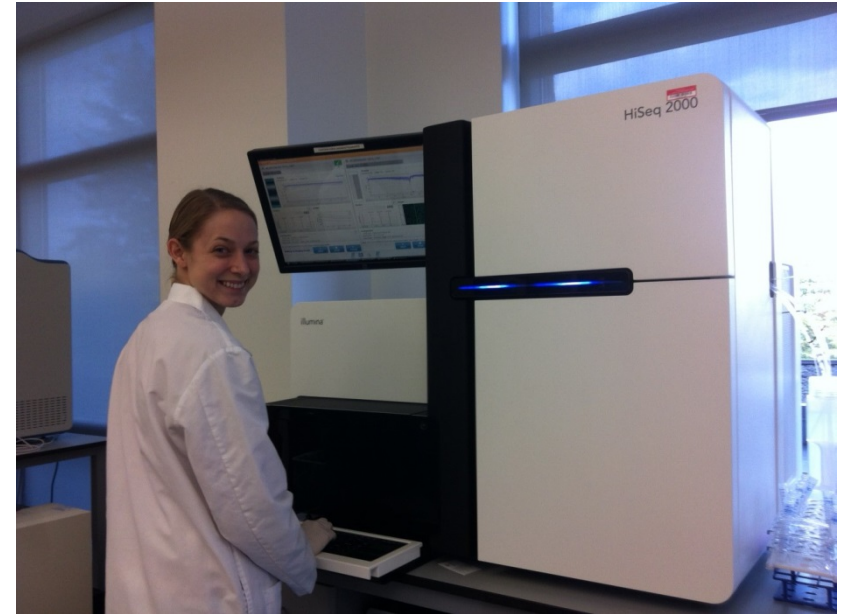## Your Body Has 10 Times As Many Microbe Cells As Human Cells

**Human genome**
23,000 genes

**Human microbiome**
1,000,000+ genes

### 99% of Your DNA Genes Are in Microbe Cells Not Human Cells

## Inclusion of the Microbiome Will Radically Change Medicine

J. Craig Venter™
INSTITUTE

it²

# A Year of Sequencing a Healthy Gut Microbiome Daily - Remarkable Stability with Abrupt Changes



**Subject A gut**

**A**

Legend:
- Proteobacteria
- Tenericutes
- Actinobacteria
- Firmicutes
- Bacteroidetes

0    100  **Days** 200    300

Living abroad

# To Map Out the Dynamics of My Microbiome Ecology I Partnered with the J. Craig Venter Institute

- **JCVI Did Metagenomic Sequencing on Seven of My Stool Samples Over 1.5 Years**

- **Sequencing on Illumina HiSeq 2000**
  - **Generates 100bp Reads**



**Illumina HiSeq 2000 at JCVI**

- **JCVI Lab Manager, Genomic Medicine**
  - **Manolito Torralba**

- **IRB PI Karen Nelson**
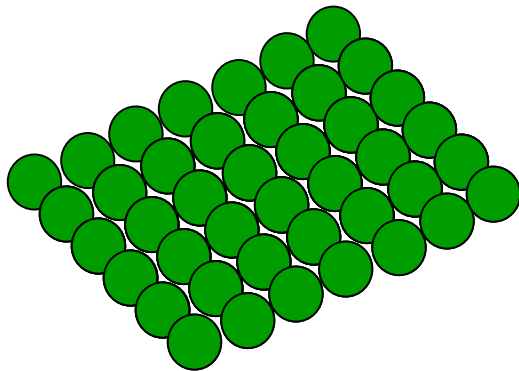  - **President JCVI**



**Manolito Torralba, JCVI**



**Karen Nelson, JCVI**

J. Craig Venter™ INSTITUTE

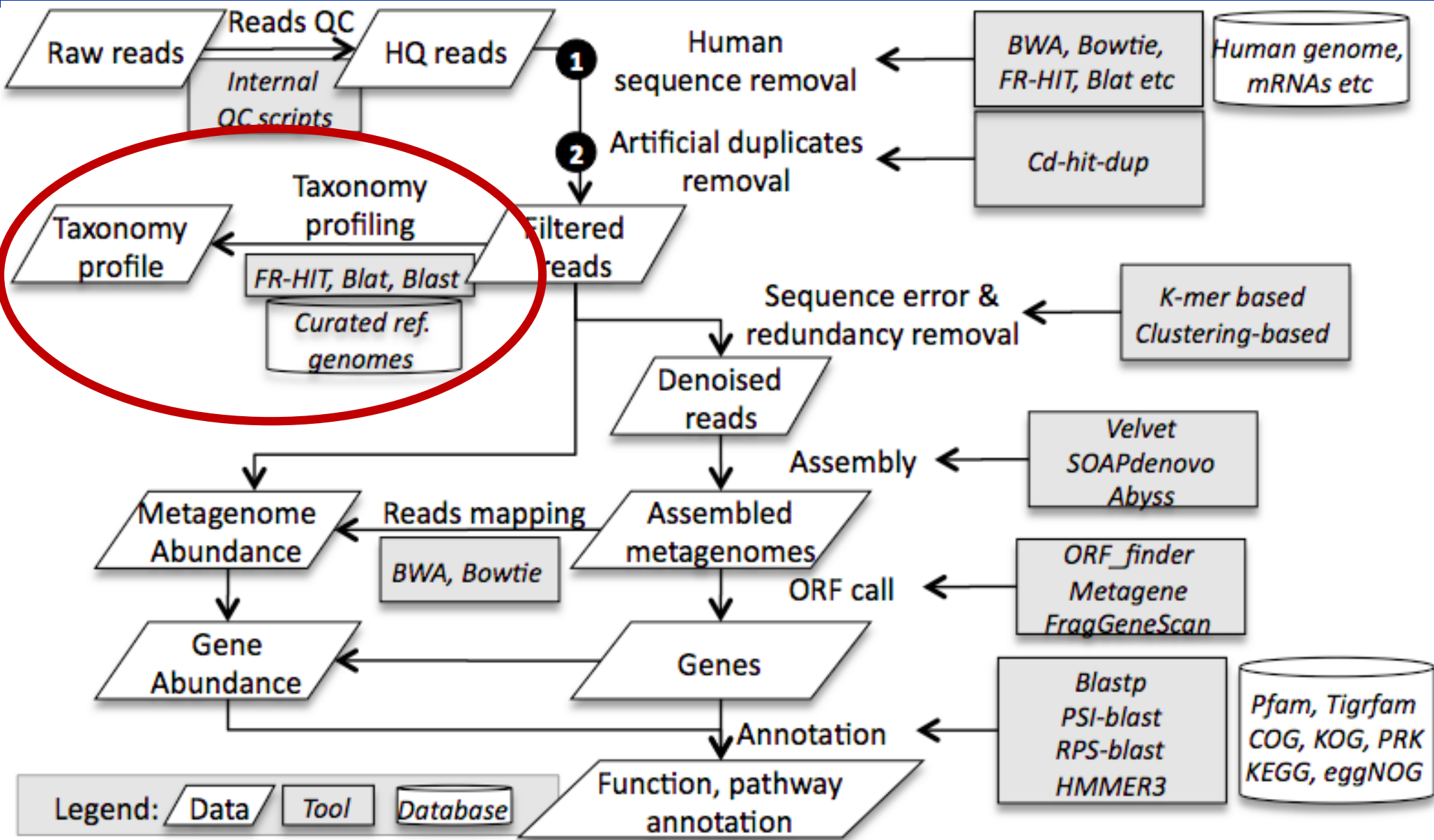# We Created a Reference Database Of Known Gut Genomes

- **NCBI April 2013**
  - **2471 Complete + 5543 Draft Bacteria & Archaea Genomes**
  - **2399 Complete Virus Genomes**
  - **26 Complete Fungi Genomes**
  - **309 HMP Eukaryote Reference Genomes**
- **Total 10,741 genomes, ~30 GB of sequences**

**Now to Align Our 5 Billion Reads Against the Reference Database**

NCBI
National Center for Biotechnology Information

HMP

NIH

J. Craig Venter™
INSTITUTE

it²

# Computational NextGen Sequencing Pipeline: From "Big Equations" to "Big Data" Computing



PI: (Weizhong Li, CRBS, UCSD):
NIH R01HG005978 (2010-2013, $1.1M)

# We Used SDSC's Gordon Data-Intensive Supercomputer to Analyze a Wide Range of Gut Microbiomes

- **~180,000 Core-Hrs on Gordon**
  - **KEGG function annotation: 90,000 hrs**
  - **Mapping: 36,000 hrs**
    - **Used 16 Cores/Node and up to 50 nodes**
  - **Duplicates removal: 18,000 hrs**
  - **Assembly: 18,000 hrs**
  - **Other: 18,000 hrs**
- **Gordon RAM Required**
  - **64GB RAM for Reference DB**
  - **192GB RAM for Assembly**
- **Gordon Disk Required**
  - **Ultra-Fast Disk Holds Ref DB for All Nodes**
  - **8TB for All Subjects**

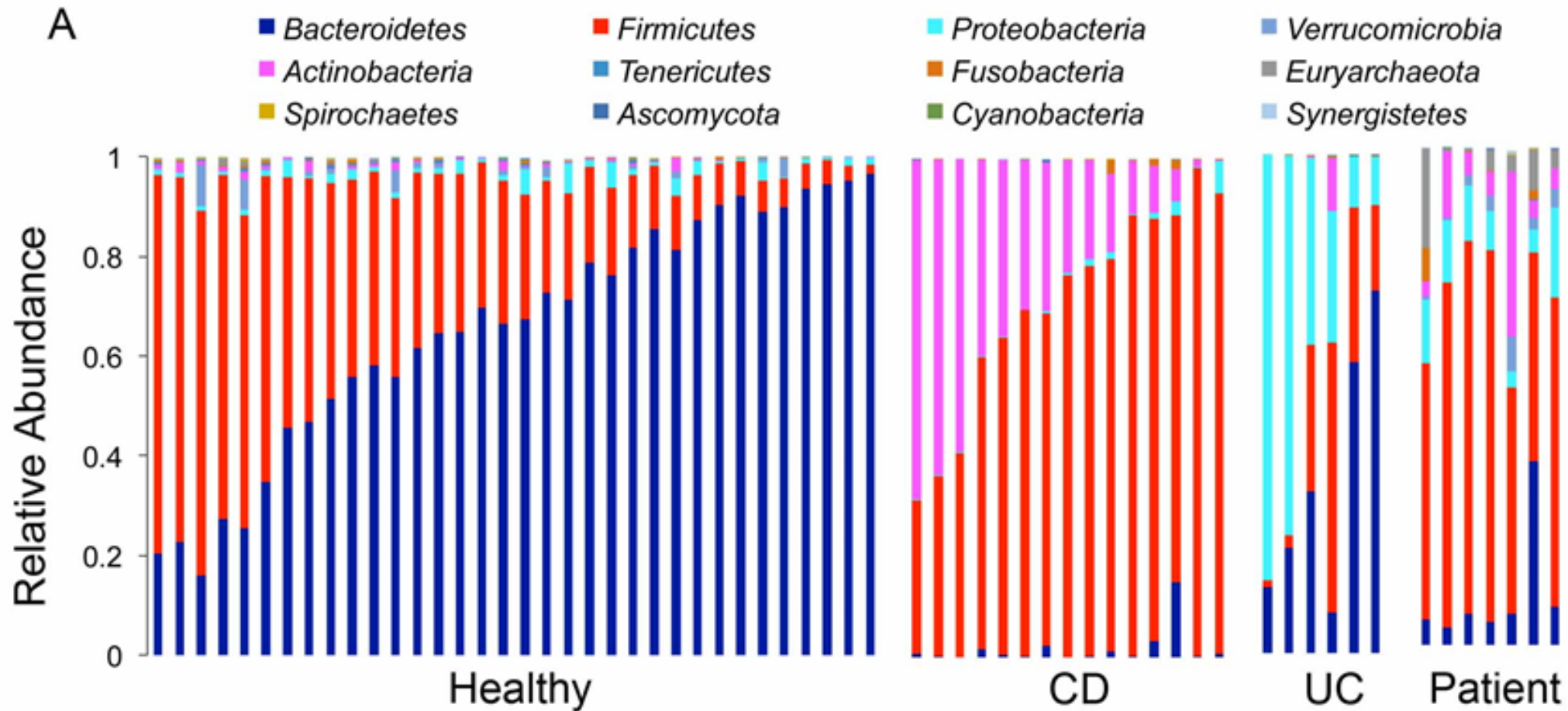**Enabled by a Grant of Time on Gordon from SDSC Director Mike Norman**

# The Emergence of Microbial Genomics Diagnostics



Microbial Ecology Is Radically Altered in Disease States, But Differently in the Two Forms of IBD
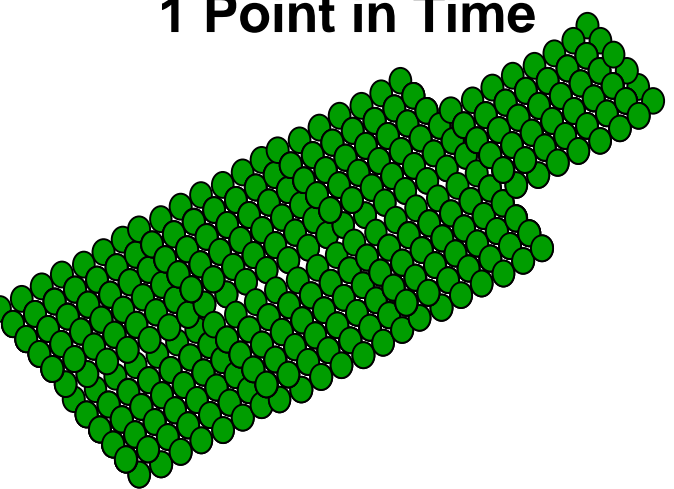
Source: Chang, et al. (2014)

# We Expaned Our Healthy Cohort to All Gut Microbiomes from NIH HMP For Comparative Analysis

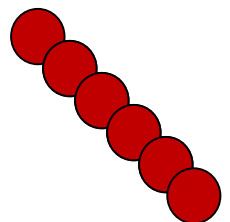**Each Sample Has 100-200 Million Illumina Short Reads (100 bases)**
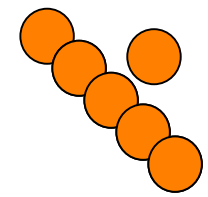
**"Healthy" Individuals**

**250 Subjects**
**1 Point in Time**

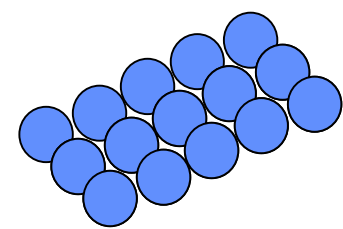**Larry Smarr**

**7 Points in Time**

**IBD Patients**

**2 Ulcerative Colitis Patients,**
**6 Points in Time**

**5 Ileal Crohn's Patients,**
**3 Points in Time**

**Total of 27 Billion Reads**
**Or 2.7 Trillion Bases**

NIH

HMP

it²

# We Used Dell's HPC Cloud to Analyze All of Our Human Gut Microbiomes

- **Dell's Sanger Cluster**
  - **32 Nodes, 512 Cores**
  - **48GB RAM per Node**

- **We Processed the Taxonomic Relative Abundance**
  - **Used ~35,000 Core-Hours on Dell's Sanger**

- **Produced Relative Abundance of ~10,000 Bacteria, Archaea, Viruses in ~300 People**
  - **~3Million Spreadsheet Cells**

- **New System: R Bio-Gen System**
  - **48 Nodes, 768 Cores**
  - **128 GB RAM per Node**
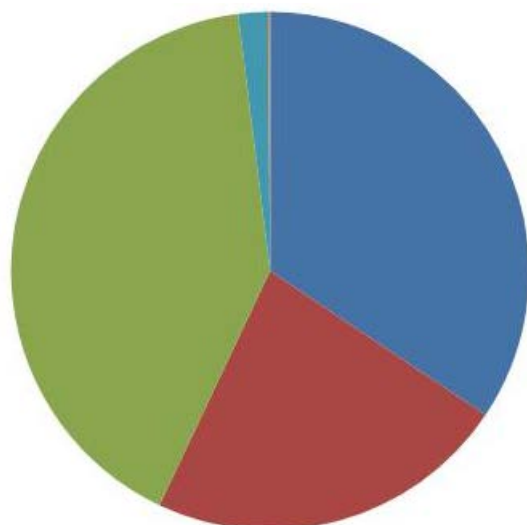
We Found Major State Shifts in Microbial Ecology Phyla Between Healthy and Two Forms of IBD
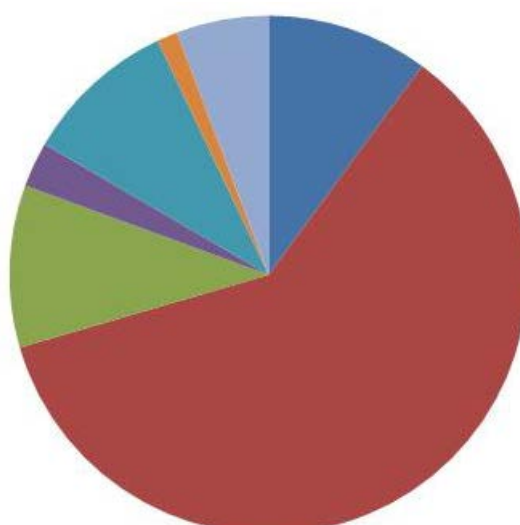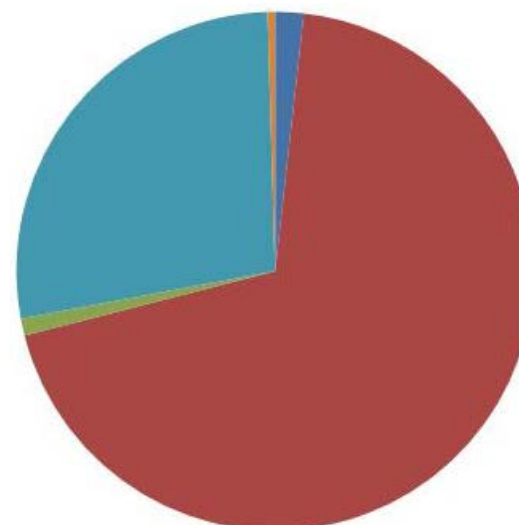
Average HE

Most Common Microbial Phyla

Legend:
- Bacteroidetes
- Firmicutes
- Proteobacteria
- Verrucomicrobia
- Actinobacteria
- Fusobacteria
- Euryarchaeota

Average Ulcerative Colitis

Average LS

Average Crohn's Disease

Explosion of Proteobacteria

Hybrid of UC and CD High Level of Archaea
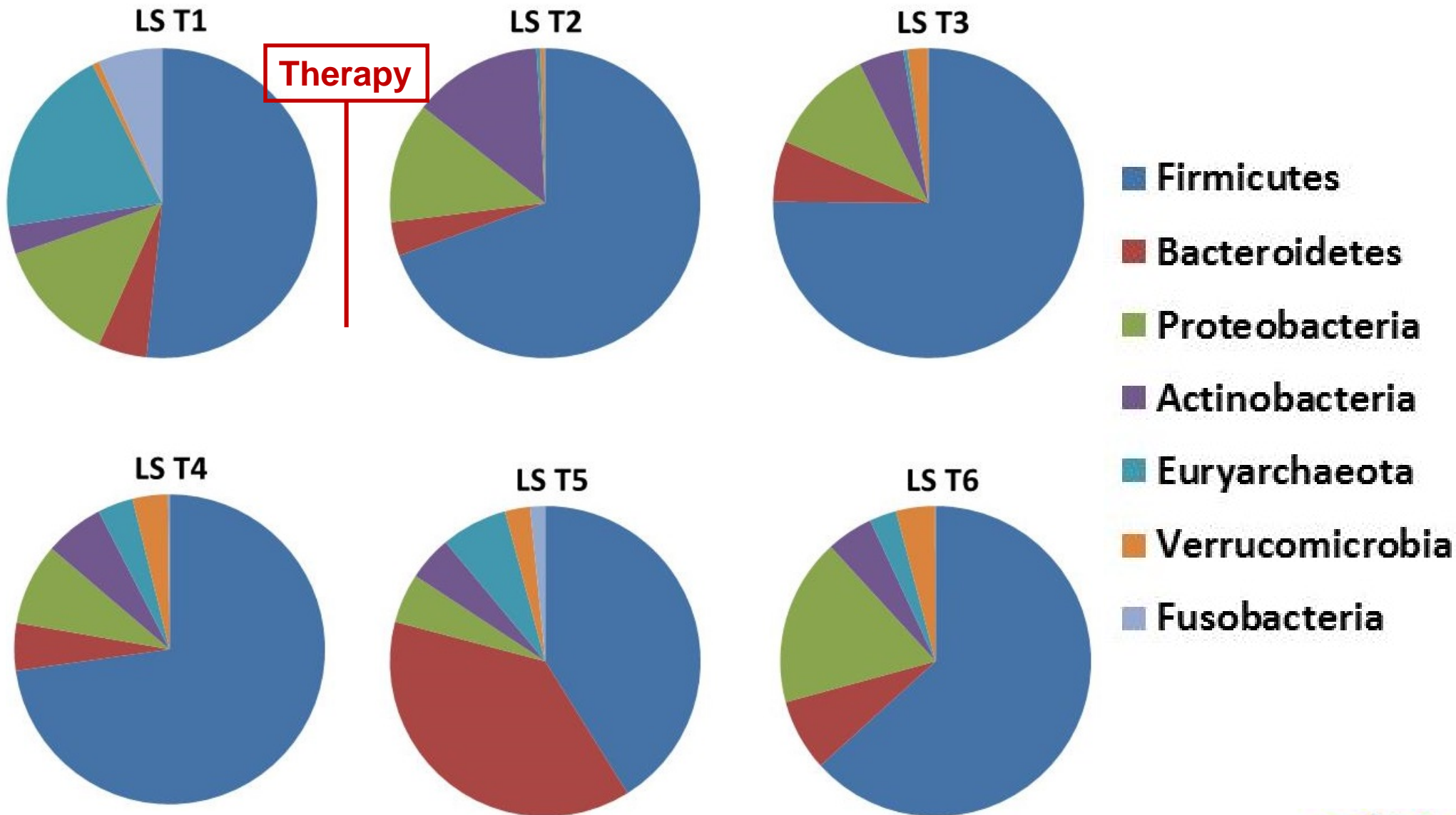
Collapse of Bacteroidetes Explosion of Actinobacteria

# Time Series Reveals Autoimmune Dynamics of Gut Microbiome by Phyla
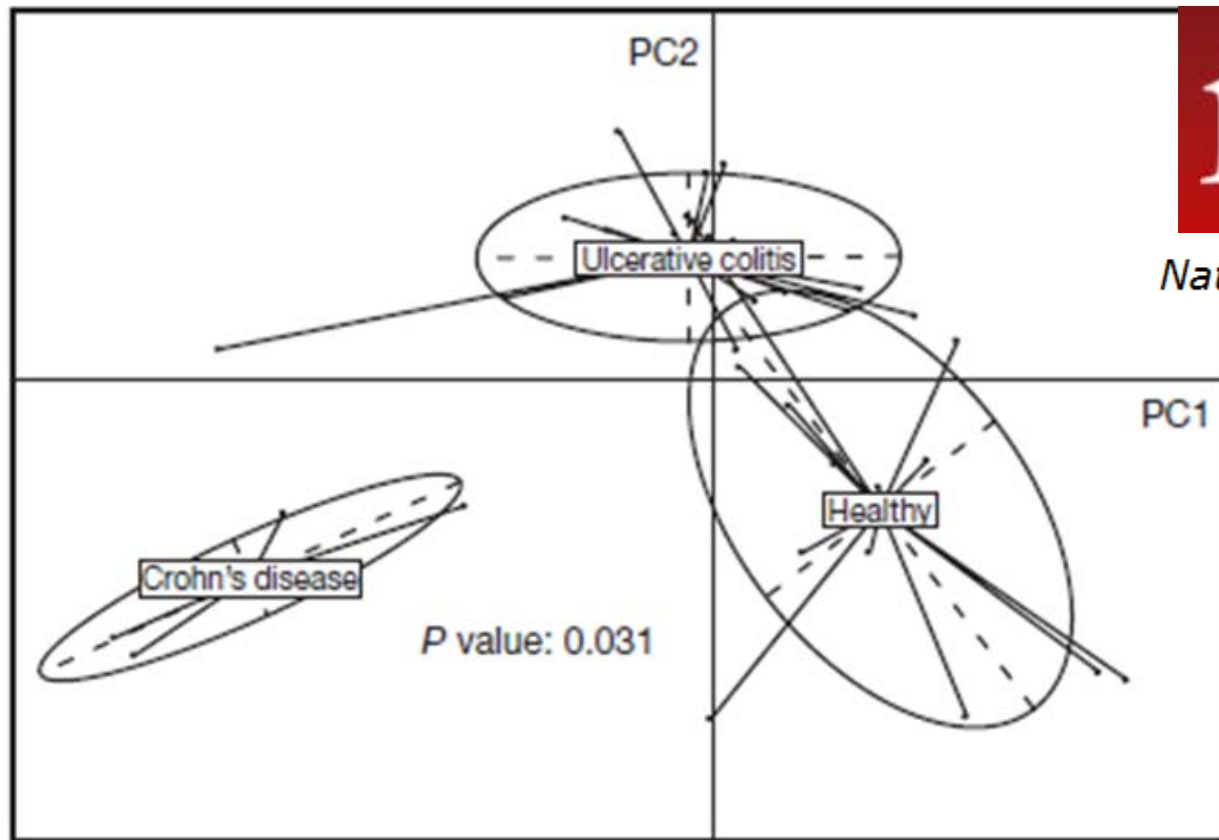
**Six Metagenomic Time Samples Over 16 Months**

**Comparing 3 LS Time Snapshots (Left)
with Healthy, Crohn's, UC (Right Top to Bottom)**

**Calit2 VROOM-FuturePatient Expedition**

Figure 4 | Bacterial species abundance differentiates IBD patients and healthy individuals. Principal component analysis with health status as

# A human gut microbial gene catalogue established by metagenomic sequencing

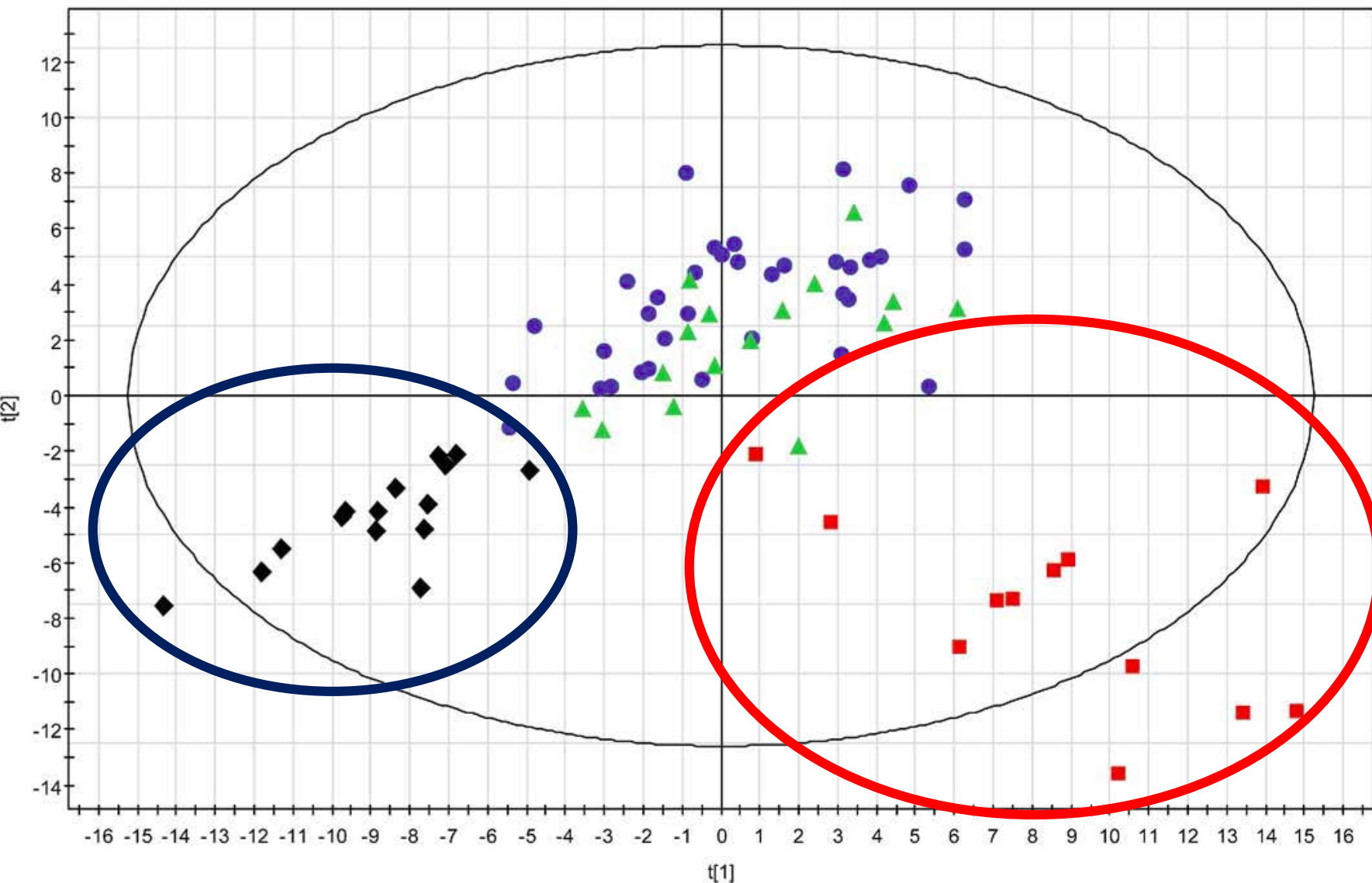Junjie Qin[1]*, Ruiqiang Li[1]*, Jeroen Raes[2,3], Manimozhiyan Arumugam[2], Kristoffer Solvsten Burgdorf[4],

# Is the Gut Microbial Ecology Different in Crohn's Disease Subtypes?

(PLS-DA)
t[Comp. 1]/t[Comp. 2]

Legend:
- ■ CCD
- ● H
- ◆ ICD
- ▲ UC

# PCA Analysis
## on Species Abundance Across People



**Green-Healthy
Red-CD
Purple-UC
Blue-LS**

ICD

CCD

Healthy Subset?

PCA2

PCA1

**Analysis by Mehrdad Yazdani, Calit2**

# Finding Species Which Differentiate Subsets of Healthy and Disease

Green-Healthy
Red-CD
Purple-UC
Blue-LS

# Dell Cloud Results Are Leading Toward Microbiome Disease Diagnosis
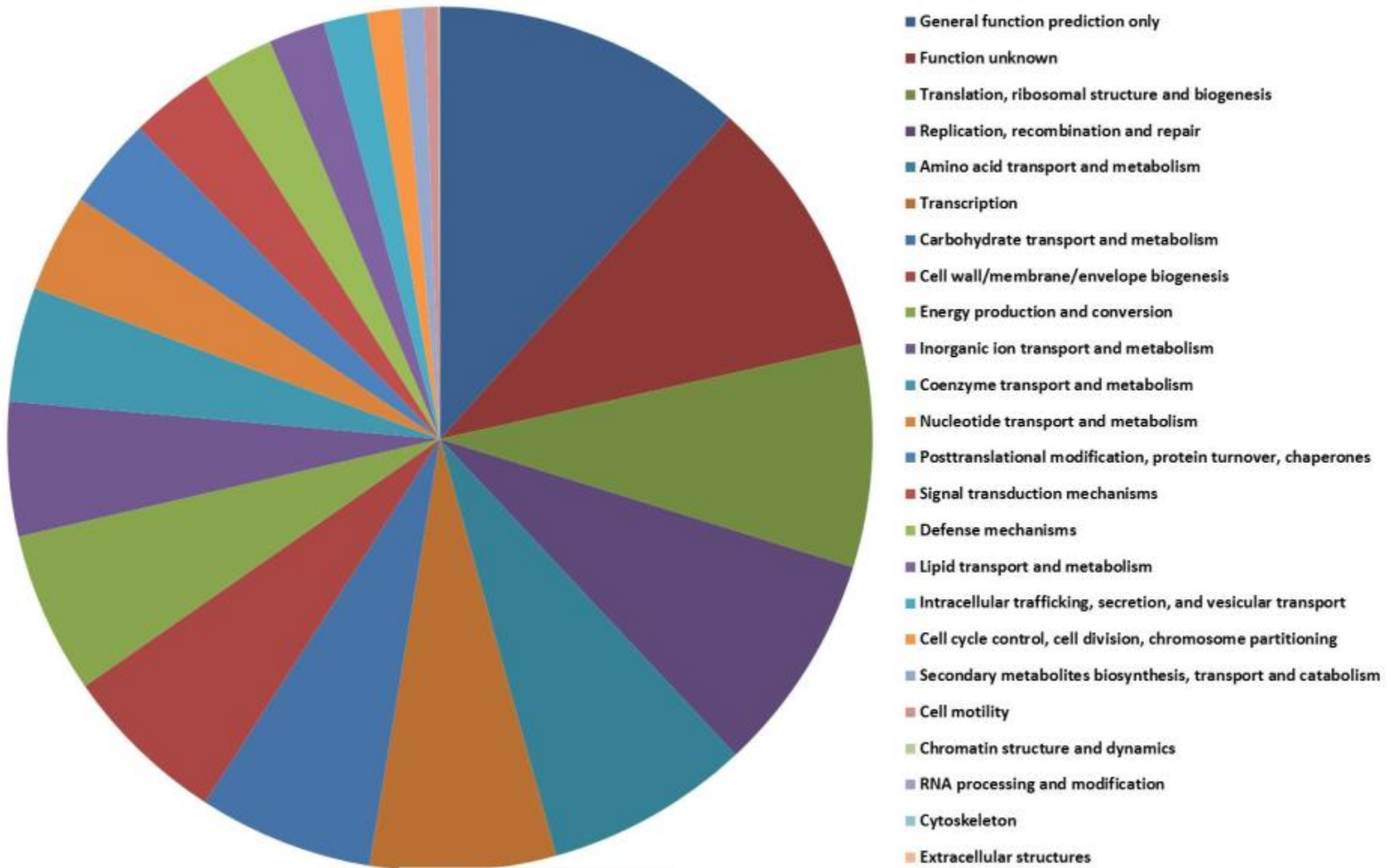


We Produced Similar Results for ~2500 Microbial Species

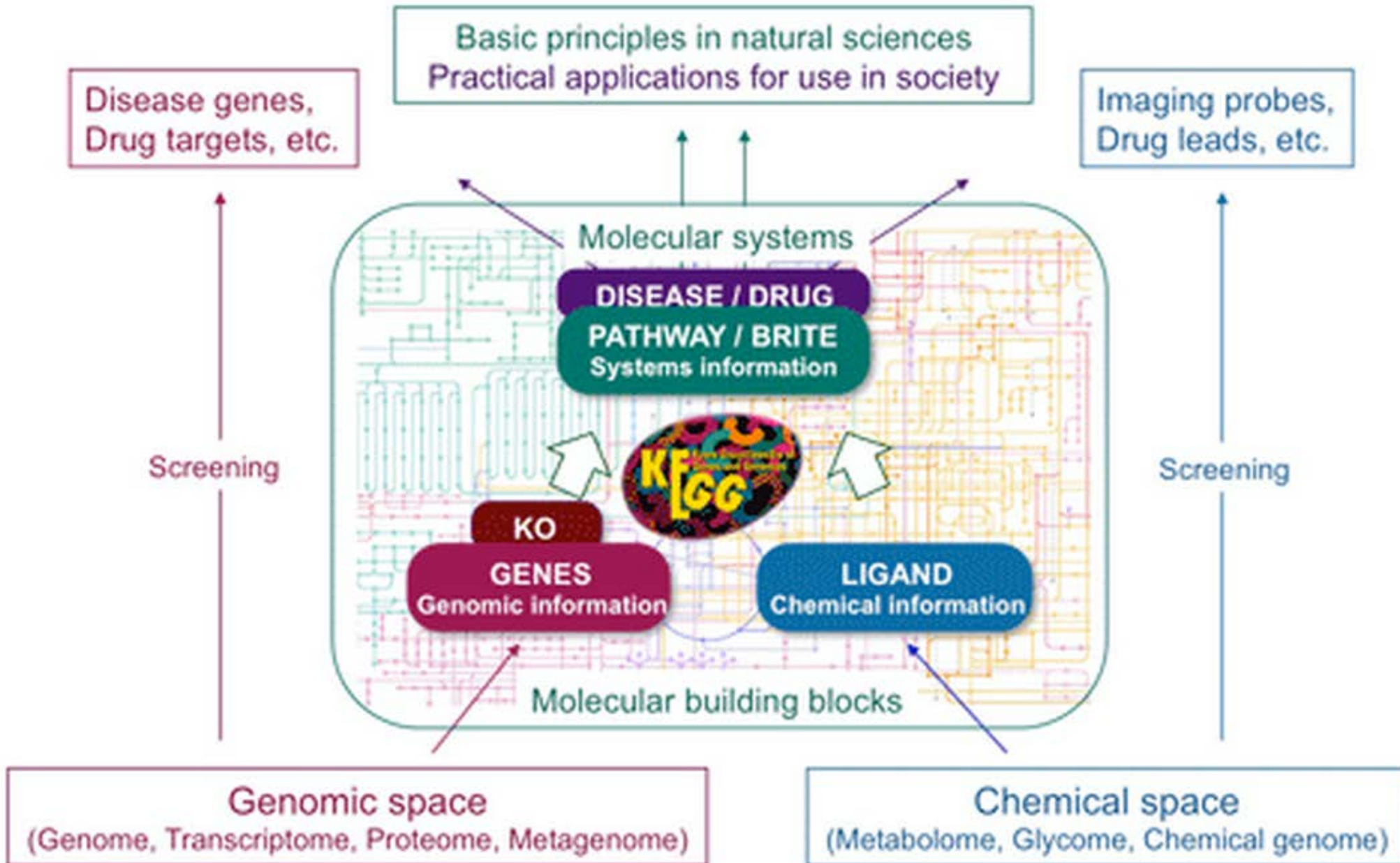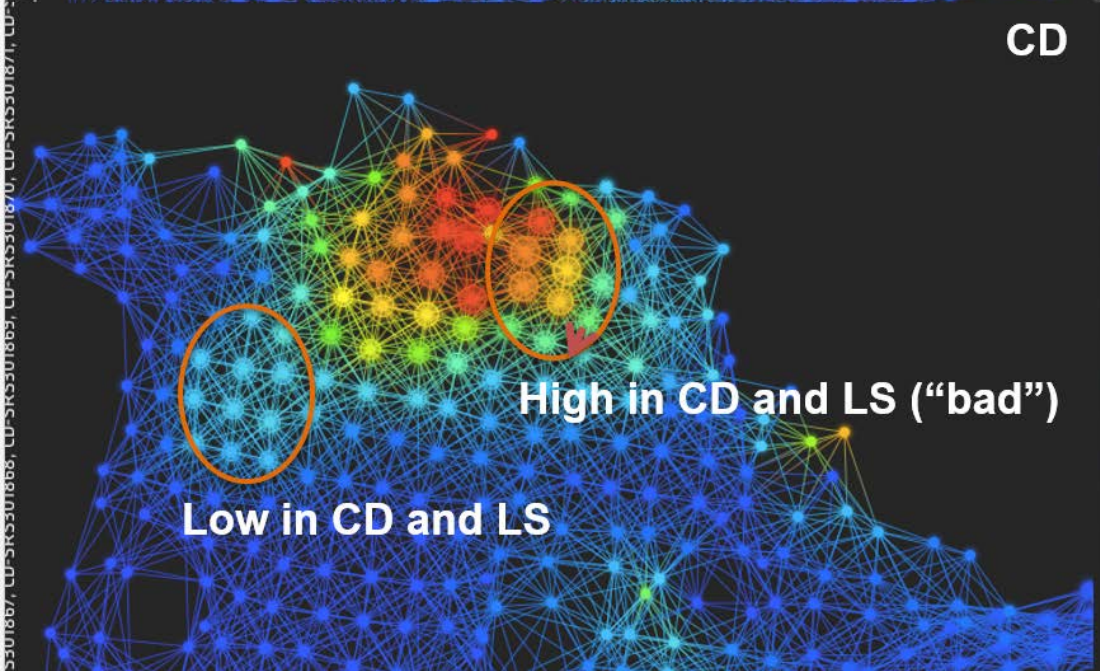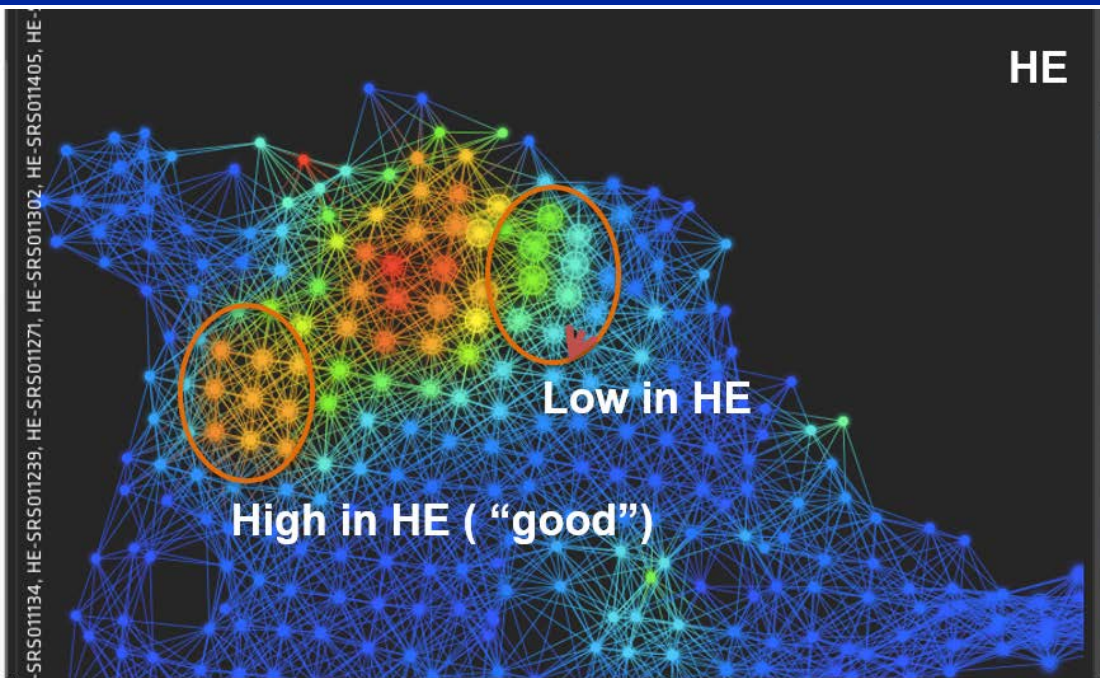# From Taxonomy to Function: Analysis of LS Clusters of Orthologous Groups (COGs)



- General function prediction only
- Function unknown
- Translation, ribosomal structure and biogenesis
- Replication, recombination and repair
- Amino acid transport and metabolism
- Transcription
- Carbohydrate transport and metabolism
- Cell wall/membrane/envelope biogenesis
- Energy production and conversion
- Inorganic ion transport and metabolism
- Coenzyme transport and metabolism
- Nucleotide transport and metabolism
- Posttranslational modification, protein turnover, chaperones
- Signal transduction mechanisms
- Defense mechanisms
- Lipid transport and metabolism
- Intracellular trafficking, secretion, and vesicular transport
- Cell cycle control, cell division, chromosome partitioning
- Secondary metabolites biosynthesis, transport and catabolism
- Cell motility
- Chromatin structure and dynamics
- RNA processing and modification
- Cytoskeleton
- Extracellular structures

J. Craig Venter INSTITUTE

UC San Diego HEALTH SCIENCES

SDSC

it²

**Analysis: Weizhong Li & Sitao Wu, UCSD**

# KEGG: a Database Resource for Understanding High-Level Functions and Utilities of the Biological System

http://www.genome.jp/kegg/

# Using Ayasdi To Discover Patterns in KEGG Dataset



HE

Low in HE

High in HE ( "good")

CD

High in CD and LS ("bad")

Low in CD and LS

Ayasdi Advanced Analytics
topological data analysis

Ayasdi Cure™
Turn Data into Therapies

**Source: Pek Lum, Chief Data Scientist, Ayasdi**

UC

**Dataset from Larry Smarr Team
With 60 Subjects (HE, CD, UC, LS)
Each with 10,000 KEGGs -
600,000 Cells**

**Full Processing to Function
(COGs, KEGGs)**

**Would Require
~1-2 Million
Core-Hours**

**Plus Dedicated Network to Move Data
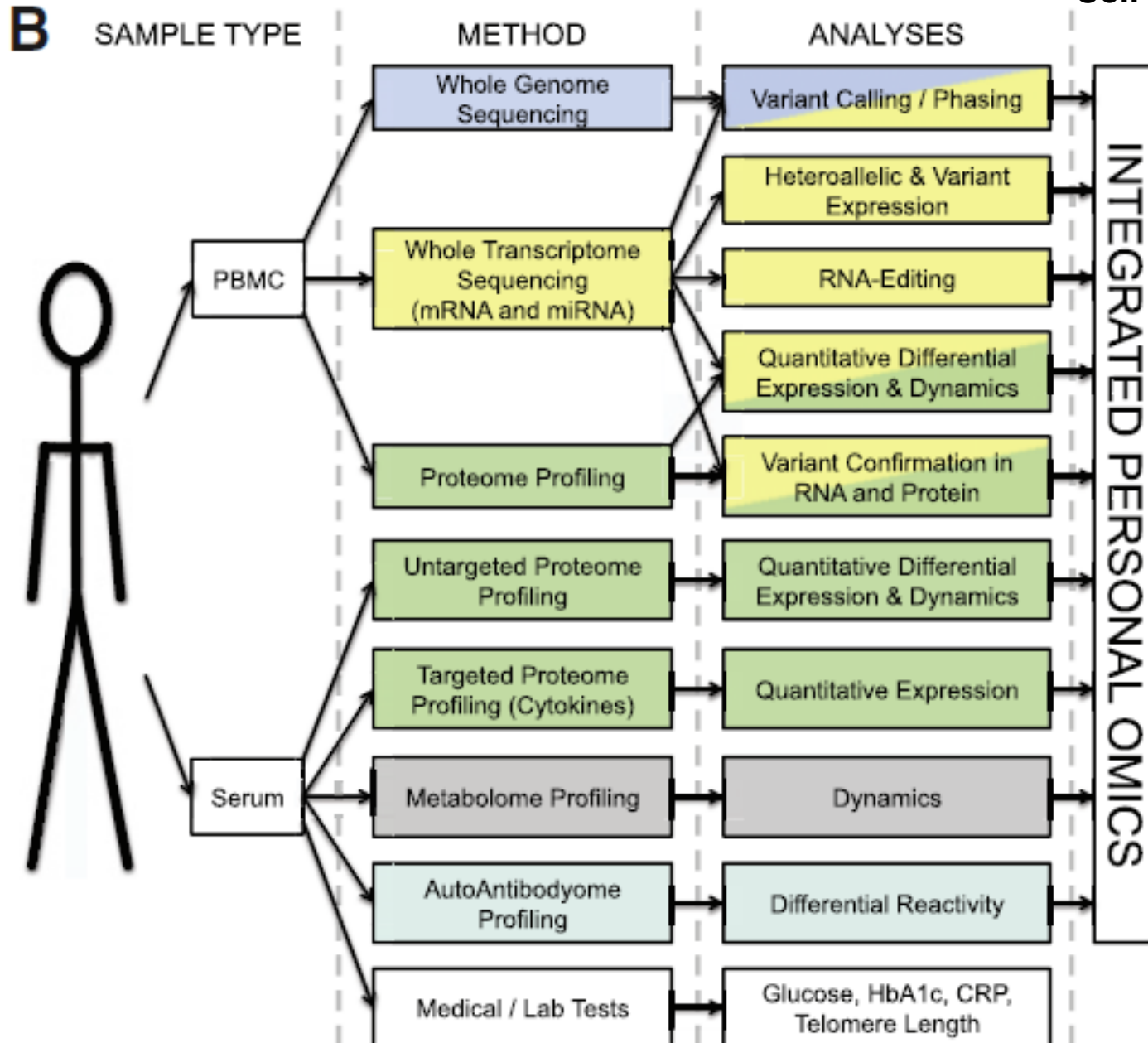From R Systems / Dell to Calit2@UC San Diego**

**Goal:  Understand
The Coupled Human Immune-Microbiome Dynamics
In the Presence of Human Genetic Predispositions**

**Drs. William J. Sandborn, John Chang, & Brigid Boland
UCSD School of Medicine, Division of Gastroenterology**

# 100x Beyond Current Medical Tests:
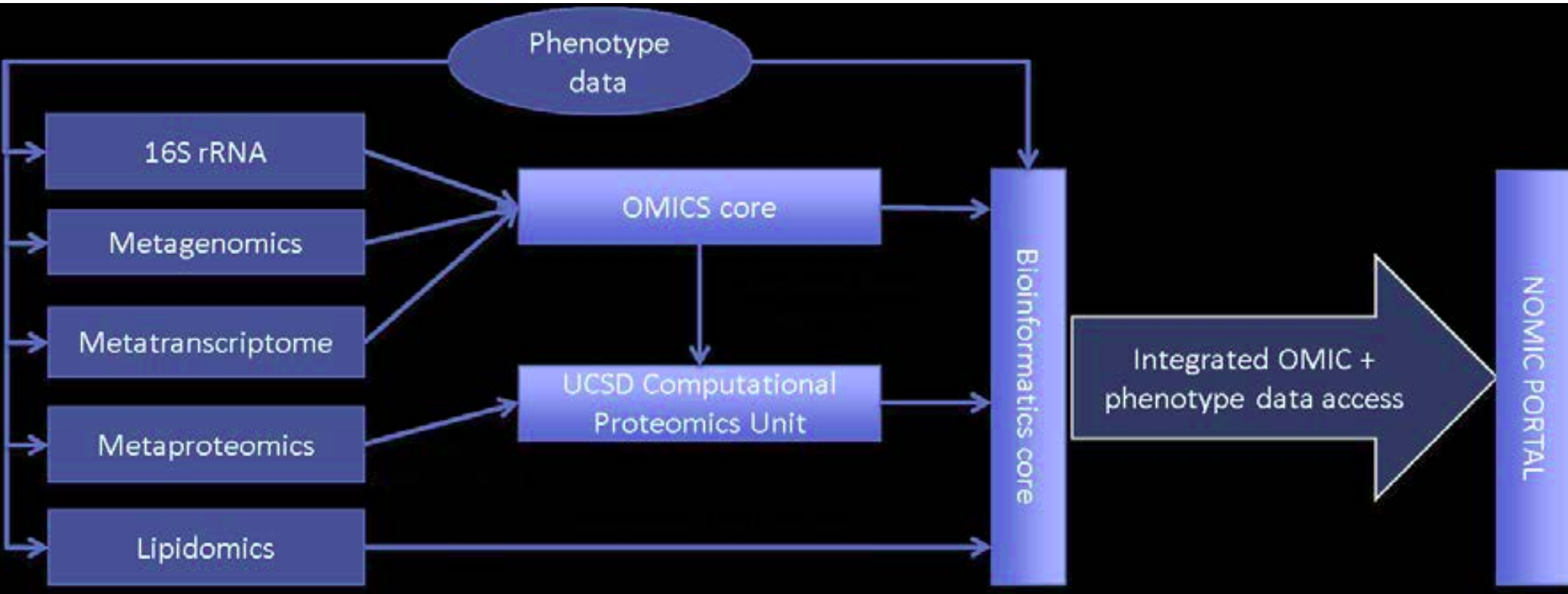# Integrated Personal Time Series of Multiple 'Omics

- **Michael Snyder, Chair of Genomics Stanford Univ.**
- **Blood Tests Time Series Over 40 Months**
  - **Tracked nearly 20,000 distinct transcripts coding for 12,000 genes**
  - **Measured the relative levels of more than 6,000 proteins and 1,000 metabolites in Snyder's blood**

# Proposed UCSD
# Integrated Omics Pipeline



**Source: Nuno Bandiera, UCSD**

# Thanks to Our Great Team!

## UCSD Metagenomics Team

**Weizhong Li**
**Sitao Wu**

## JCVI Team

**Karen Nelson**
**Shibu Yooseph**
**Manolito Torralba**

## Calit2@UCSD
## Future Patient Team

**Jerry Sheehan**
**Tom DeFanti**
**Kevin Patrick**
**Jurgen Schulze**
**Andrew Prudhomme**
**Philip Weber**
**Fred Raab**
**Joe Keefe**
**Ernesto Ramirez**

## SDSC Team

**Michael Norman**
**Mahidhar Tatineni**
**Robert Sinkovits**

## UCSD Health Sciences Team

**William J. Sandborn**
**Elisabeth Evans**
**John Chang**
**Brigid Boland**
**David Brenner**