# High-throughput sequencing and big data: implications for personalized medicine?

## Dominick J. Lemas, PhD

*Postdoctoral Fellow in Pediatrics - Neonatology*
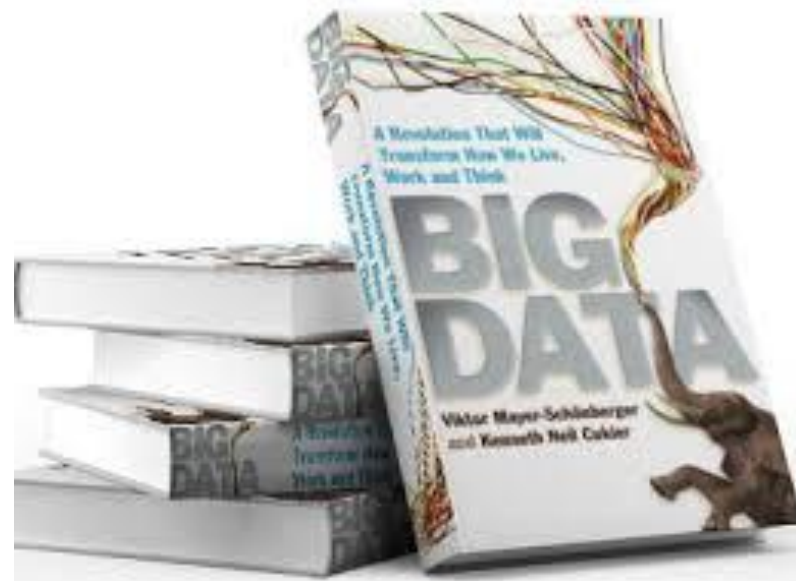
*Mentor: Jacob E. (Jed) Friedman, PhD*

University of Colorado
Anschutz Medical Campus

# What is Big Data?

1) Top Tech phrase of 2013[1]

2) Messy, Noisy, Imprecise

3) Datafication

4) Repurposed

5) N= All

6) Privacy?

[1]http://www.languagemonitor.com/
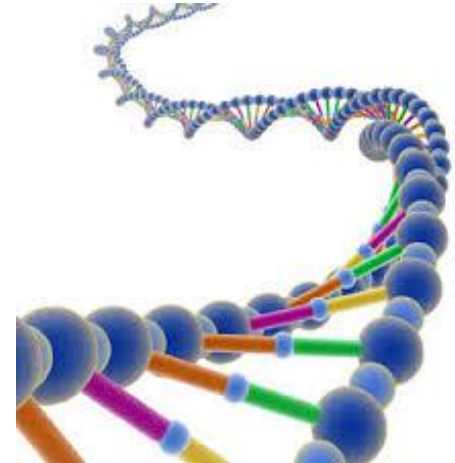
University of Colorado
Anschutz Medical Campus

# What is the human genome?

1. Human body contains trillions of cells.
2. Each cell contains a nucleus.
3. Each nucleus contains 2 complete genomes.
4. . . . . . . .

**OR  if the genome was a book!**

- **There are 23 chapters → chromosomes**
- **Each chapter contains several hundred stories→ genes**
- **Each story is composed of paragraphs → exons**
- **Interrupted by advertisements → introns**
- **Each paragraph is made up of words → codons**
- **Each word is written in letters → bases  . . . A,C,T,G**

Genome. Matt Ridley. 1999

University of Colorado
Anschutz Medical Campus

# Goals of Human Genome Project

1) Generate working draft of 90% of the human genome (2001).

2) Obtain complete, high-quality genomic sequence (2003).

3) Make all data publically available.

4) Develop novel sequencing technologies.

5) Map Sequence Variation.

6) Interpret functions of genome.

7) Develop comparative genomic strategies.

8) Ethical, legal and social implications (ELSI).
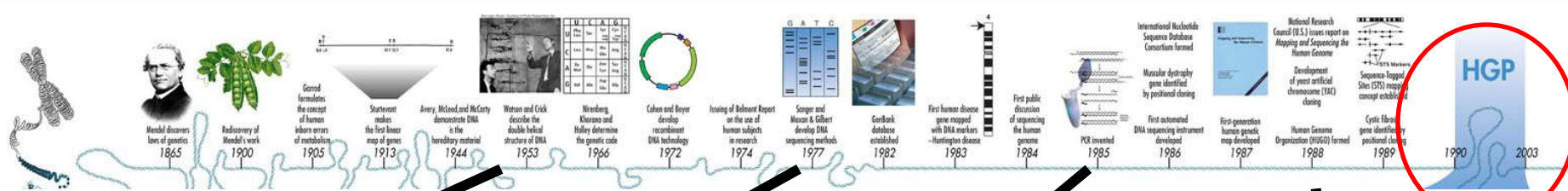
9) Bioinformatics and Computational Biology

10) Training

# Benefits of Sequencing Human Genome

1) Molecular Medicine

2) Energy and Environmental Applications

3) Bioarcheology, anthropology, evolution, human migration

4) DNA forensics

5) Agriculture, livestock breeding, and bioprocessing

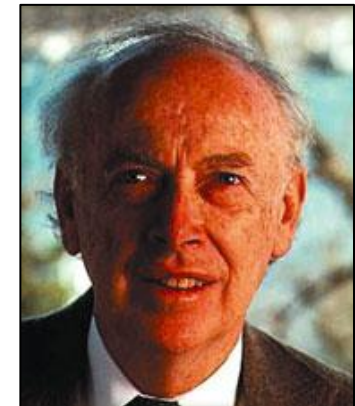# Sequencing Milestones: the early days



1953

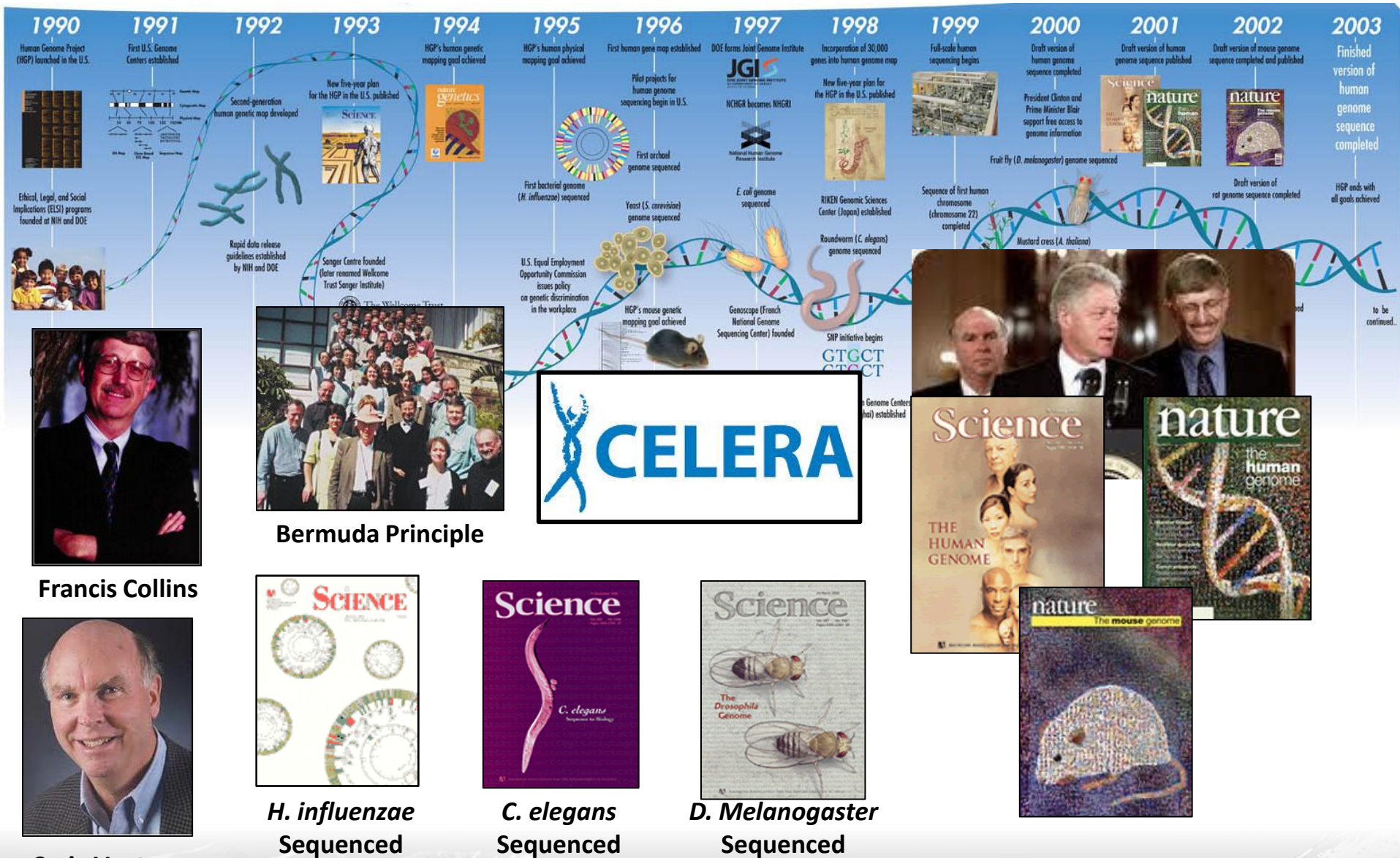Watson & Crick

1977

Fredrick Sanger

1985

Charles DeLisi

1988

James D. Watson

*Collins. 2001.Nature* **422**, 835-847

University of Colorado
Anschutz Medical Campus

# Sequencing Milestones: HGP



Timeline 1990–2003:

**1990** — Human Genome Project (HGP) launched in the U.S. Ethical, Legal, and Social Implications (ELSI) programs founded at NIH and DOE

**1991** — First U.S. Genome Centers established

**1992** — Second-generation human genetic map developed. Rapid data release guidelines established by NIH and DOE. Sanger Centre founded (later renamed Wellcome Trust Sanger Institute)

**1993** — New five-year plan for the HGP in the U.S. published

**1994** — HGP's human genetic mapping goal achieved

**1995** — HGP's human physical mapping goal achieved. Pilot projects for human genome sequencing begin in U.S. First bacterial genome (H. influenzae) sequenced. U.S. Equal Employment Opportunity Commission issues policy on genetic discrimination in the workplace. HGP's mouse genetic mapping goal achieved

**1996** — First human gene map established. First archeal genome sequenced. Yeast (S. cerevisiae) genome sequenced

**1997** — DOE forms Joint Genome Institute (JGI). NCHGR becomes NHGRI National Human Genome Research Institute. E. coli genome sequenced. Genoscope (French National Genome Sequencing Center) founded

**1998** — Incorporation of 30,000 genes into human genome map. New five-year plan for the HGP in the U.S. published. RIKEN Genomic Sciences Center (Japan) established. Roundworm (C. elegans) genome sequenced. SNP initiative begins

**1999** — Full-scale human sequencing begins. Sequence of first human chromosome (chromosome 22) completed

**2000** — Draft version of human genome sequence completed. President Clinton and Prime Minister Blair support free access to genome information. Fruit fly (D. melanogaster) genome sequenced. Mustard cress (A. thaliana)

**2001** — Draft version of human genome sequence published

**2002** — Draft version of mouse genome sequence completed and published. Draft version of rat genome sequence completed

**2003** — Finished version of human genome sequence completed. HGP ends with all goals achieved. to be continued...

**Francis Collins**

**Craig Venter**

**Bermuda Principle**

CELERA

**H. influenzae Sequenced**

**C. elegans Sequenced**

**D. Melanogaster Sequenced**

*Collins. 2001. Nature* **422**, 835-847

# The International Human Genome Sequencing Consortium

**G5- Completed Bulk of Sequencing**
- **Whitehead Institute/MIT Center for Genome Research, Cambridge, Mass., U.S.**
- **The Wellcome Trust Sanger Institute, Hinxton, Cambridgeshire, U.K.**
- **Washington University School of Medicine Genome Sequencing Center, St. Louis, Mo., U.S.**
- **U. S. Department of Energy Joint Genome Institute, Walnut Creek, Calif., U.S.**
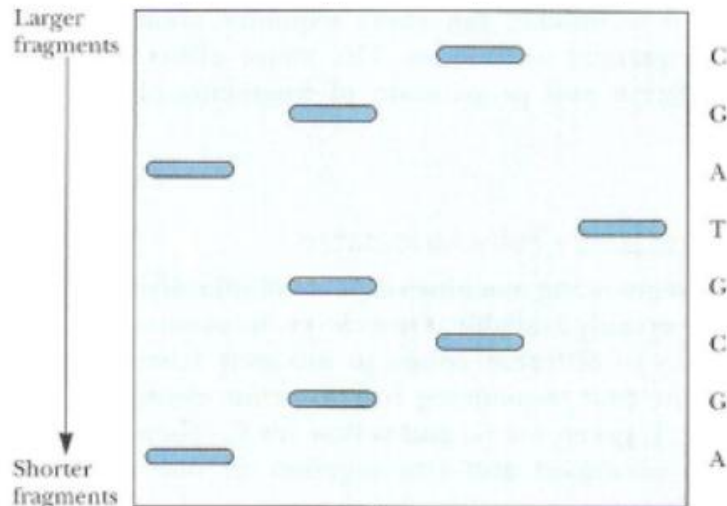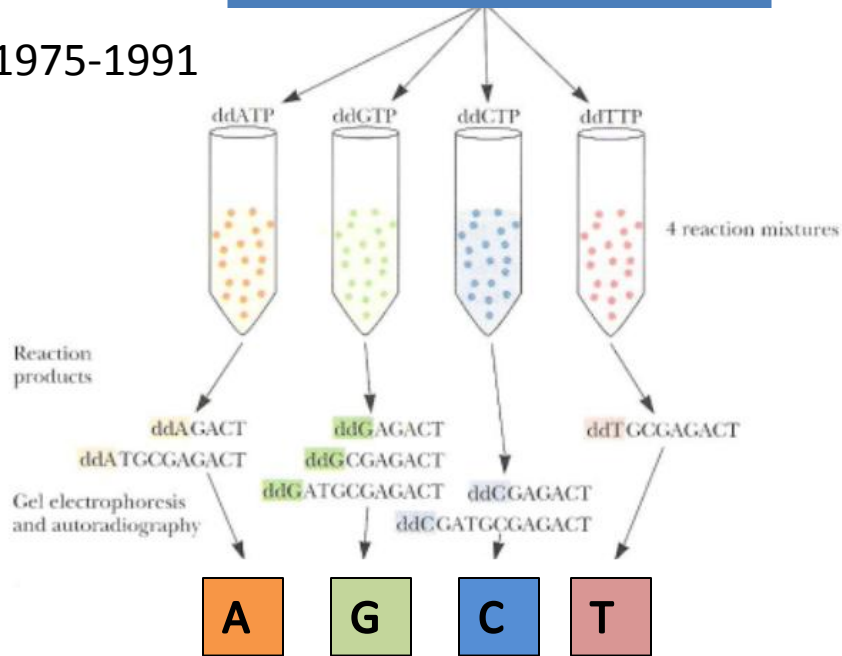- **Baylor College of Medicine Human Genome Sequencing Center, Houston, Tex., U.S.**

- RIKEN Genomic Sciences Center, Yokohama, Japan
- Genoscope and CNRS UMR-8030, Evry, France
- GTC Sequencing Center, Waltham, Mass., U.S.
- Department of Genome Analysis, Jena, Germany
- Beijing Genomics Institute/Human Genome Center, Beijing, China
- Multimegabase Sequencing Center, Seattle, Wash., U.S.
- Stanford Genome Technology Center, Stanford, Calif., U.S.
- Stanford Human Genome Center, Stanford, Calif., U.S.
- University of Washington Genome Center, Seattle, Wash., U.S.
- Department of Molecular Biology, Tokyo, Japan
- University of Texas Southwestern Medical Center at Dallas, Dallas, Texas, U.S.
- University of Oklahoma's Advanced Center for Genome Technology, Norman, Okla., U.S.
- Max Planck Institute for Molecular Genetics, Berlin, Germany
- Cold Spring Harbor Laboratory, Lita Annenberg Hazen Genome Center, Cold Spring Harbor, N.Y., U.S.
- GBF - German Research Centre for Biotechnology, Braunschweig, Germany

University of Colorado
Anschutz Medical Campus

# Sanger Sequencing



**DNA Template**

1975-1991

**The ABI Prism 3700 /3730**

1992-2000's

- $300,000/machine
- Sequence 50-100K bp/hr

University of Colorado
Anschutz Medical Campus

# Whole-Genome Sequencing

ACGTCCTATGCGTATGCGTAATGCCACATATTGCTATGCGTAATGCGTACC

genome

**Break genome into small pieces called "reads"**

*N* reads

TATGCGTATGCGTAATG

read length *L*

**Sequence reads**

**Assemble Reads**

University of Colorado
Anschutz Medical Campus

# Computational Challenges

- Coverage?

- Imputation?

- Alignment?

- Formatting?

- Analysis?

# Where does the data live?

# Sequencing has gotten Cheaper and Faster

Cost of one human genome

- HGP: $ 3 billion
- 2004: $30,000,000
- 2008: $100,000
- 2010: $10,000
- **2011: $4,000**
- 2012-13: $1,000
- ???: $300



US $100 MILLION

MOORE'S LAW

$10 MILLION

$1 MILLION

COST PER GENOME

$100 000

$10 000

$1000

2001                                              2012

# BIG DATA & Sequence



Data Size + Computing Speed

Growth of Genbank Data (doubles every 9 months)

Time to compute all Genbank data with fastest available computer (doubles every 18 months)

Growth of computer speed (Moore's Law) (doubles every 18 months)

Time

**Growth of GenBank and WGS**

Whole-Genome Shotgun (WGS)

GenBank

University of Colorado
Anschutz Medical Campus

# So What Did We Learn?

- <3% of genome encodes ~20,000 genes.

- More than half of genome is repetitive.

- Identification of ~2,850 gene impact rare diseases.

- ~1,100 markers affecting common disease & ~150 targets for cancer.

- "Big Science" can win.

- Cost of sequencing per base has been reduced by magnitude of ~100k.

*Lander. 2011. Nature* **470** *, 187-197*

# Mapping Genetic Variation



~7 billion people



~8.7 million species



~37 trillion cells/human

University of Colorado
Anschutz Medical Campus

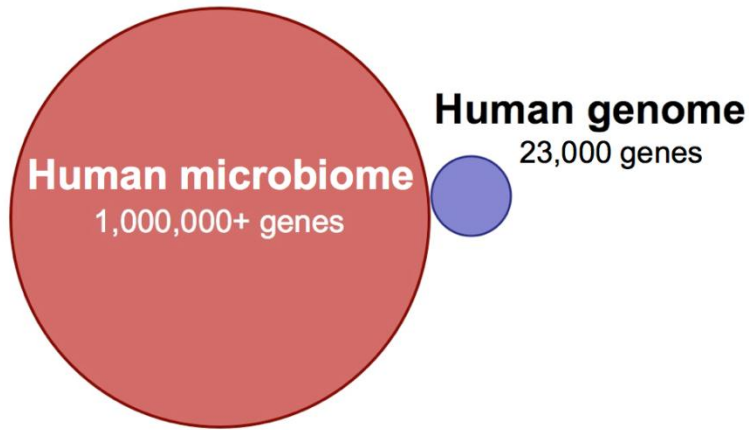# Obesity is Complex

# Gene-Environment Interactions

# Variation in the Human Microbiome

**Microbe:** tiny living organism, such as bacterium, fungus, protozoan, or virus.

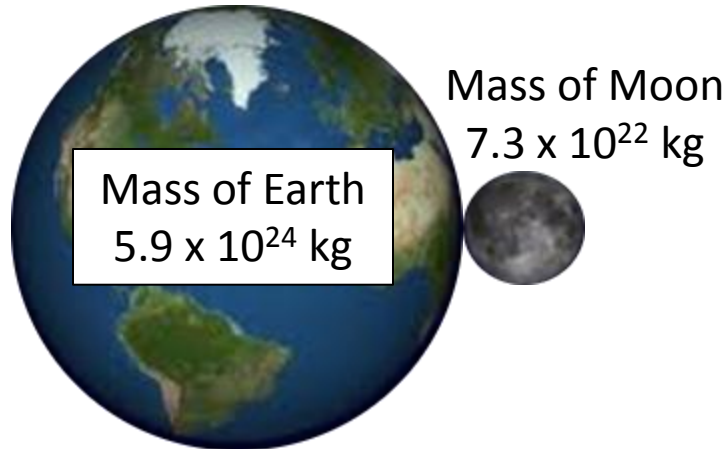**Microbiome:** collectively all the microbes in the human body; a community of microbes.

# Variation in Human Microbiome?



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

**Human genome**
23,000 genes

Mass of Moon
$7.3 \times 10^{22}$ kg

Mass of Earth
$5.9 \times 10^{24}$ kg

**Relative to humans:**

- 9 in 10 cells are microbial!
- ~1000 different species.
- ~150x more genes.
- ~3lbs of microbes in the human gut.
- ~60% of stool by dry mass is microbial.

**Primary functions of the microbiome:**
- Stimulate the development of our immune system.
- Resist colonization by pathogens.
- Extract energy from food.

University of Colorado
Anschutz Medical Campus

# Does the microbiome impact maternal metabolism during pregnancy ?



Koren et al 2012

**We hypothesize the maternal microbiome in mothers will directly affect the development of the infants microbiome and adiposity during the first year of life.**
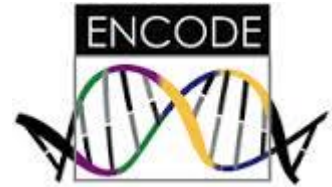
# Overview of Specific Aims

# How Do You Measure the Microbiome?



Target Highly Conserved Bacterial Gene
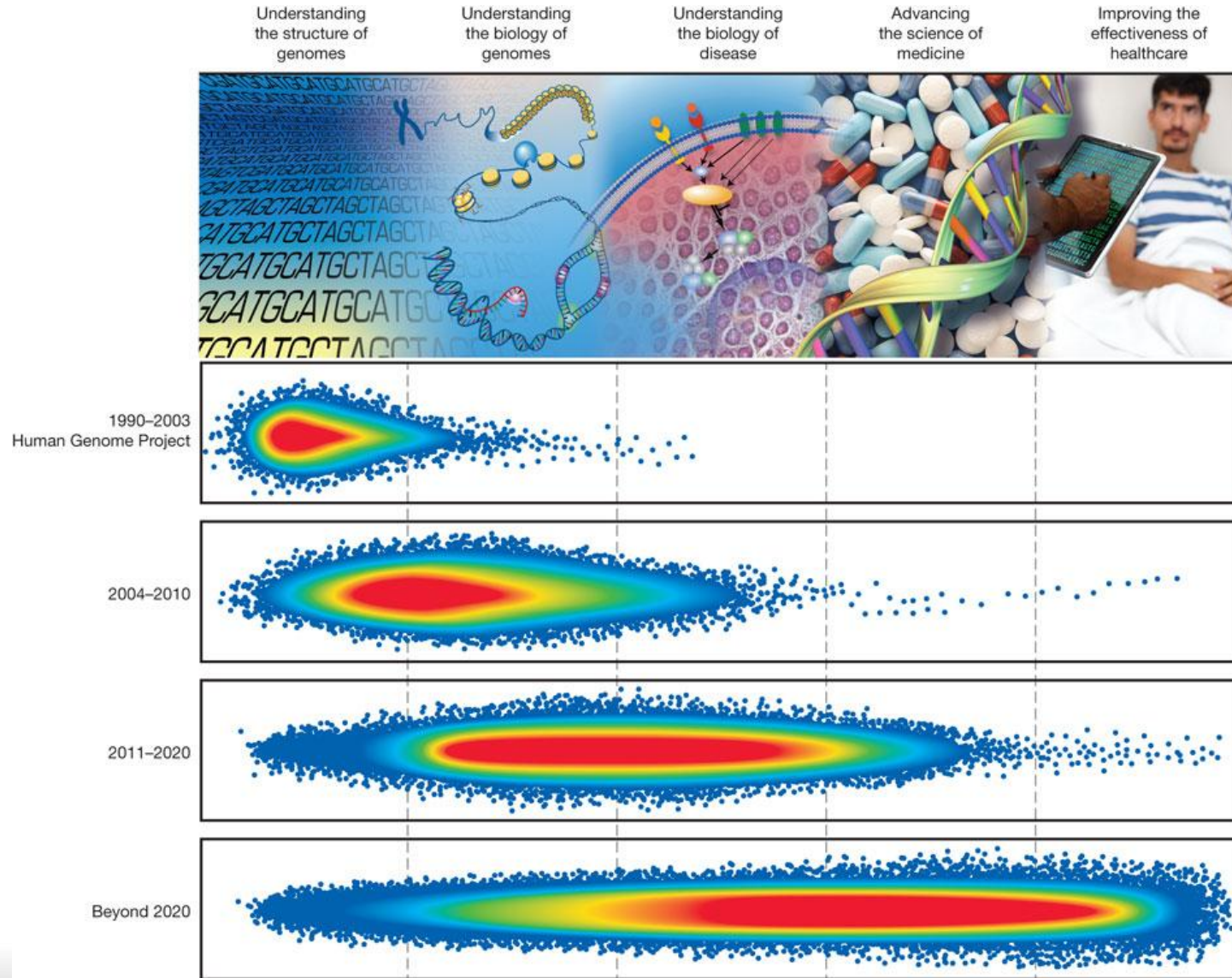
DNA samples

Target All Bacterial Genomes

Illumina MiSeq Personal Sequencer

Who is present?

What are they doing?

University of Colorado
Anschutz Medical Campus

# Moving toward a Personalized Medicine?



**But has sequencing the human genome improved health?**

# Five Domains of Genomic Research



Green. 2011. *Nature* **470**, 204-213

University of Colorado
Anschutz Medical Campus

# Bioinformatics and Computational Biology to the Rescue?

**Data Analysis-** existing tools are becoming inadequate to analyze data.

**Data Integration-** Need to harmonize disparate data types.

**Visualization-** Need to accommodate multi-dimensional data.

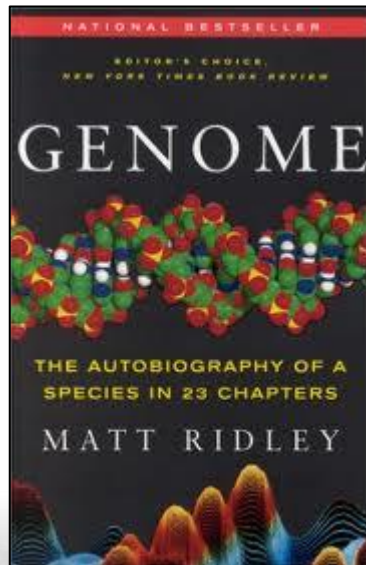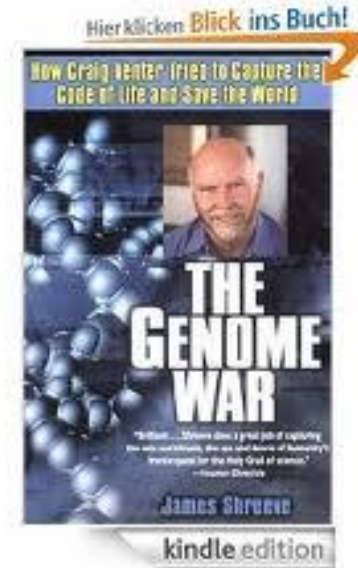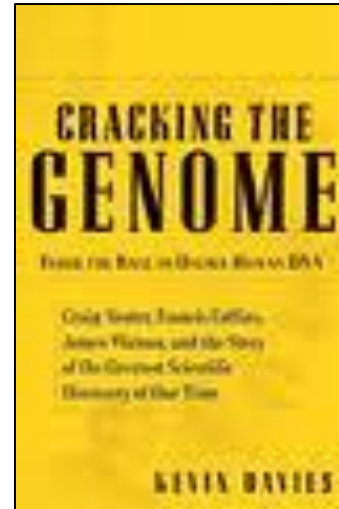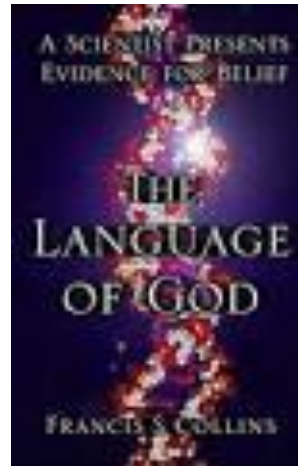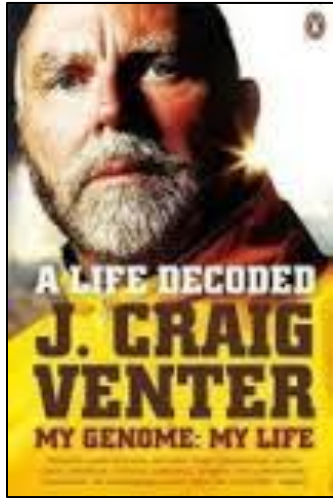**Computational Tools and Infrastructure-** storage, capacity, privacy.

**Training-** Need biologist, computer science, informatics, data science . . . . . .

Green. 2011. *Nature* **470**, 204-213

University of Colorado
Anschutz Medical Campus

# Recommended Reading

# Special Thanks to:



My mentors
- Jacob E. Friedman, PhD
- Daniel N. Frank, PhD
- Stephanie A. Santorico, PhD
- Linda A. Barbour, MD, MSPH
- Dana Dabelea, MD, PhD

Funding Support:
- ADA/Glaxo-SmithKline
- T32-HD07186



Friedman Lab 2012

# Questions?

University of Colorado
Anschutz Medical Campus