

The Genetic and Evolutionary Computation Conference (GECCO'12)

BIOINFORMATICS



William S. Bush, Ph.D.

Assistant Professor of Biomedical Informatics
Center for Human Genetics Research
Vanderbilt University, Nashville, Tennessee USA
<http://www.sigevo.org/gecco-2012/>
william.s.bush@vanderbilt.edu
gettinggeneticsdone.com
chgr.mc.vanderbilt.edu/bushlab

Copyright is held by the author/owner(s).
GECCO'12 Companion, July 7–11, 2012, Philadelphia, PA, USA.
ACM 978-1-4503-1178-6/12/07.



VANDERBILT UNIVERSITY

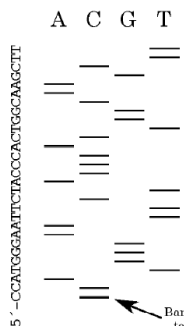


ABOUT THE INSTRUCTOR

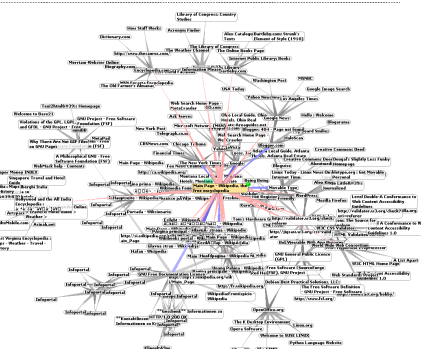
- William S. Bush is an Assistant Professor of Biomedical Informatics in the Center for Human Genetics Research at Vanderbilt University. Using a combination of bioinformatics, basic statistical approaches, and more advanced data mining and machine learning techniques, he studies how patterns of genomic variation influence the function of both individual genes and entire biological systems.



VANDERBILT UNIVERSITY



BIOINFORMATICS (1990)



VANDERBILT UNIVERSITY



Bioinformatics

Mark S Boguski

National Center for Biotechnology Information, Bethesda, USA

Computer databases, networks and software tools are essential materials and methods for biomedical research and are involved in almost every aspect of disease gene mapping and positional cloning. Public databases of DNA and protein sequences and genetic and physical map information are increasing rapidly in size and complexity and are also improving in quality, comprehensiveness, interoperability and access. A new generation of software tools for navigating through the biomedical literature has become available. Programs for sequence homology searching and genetic map construction have become more sophisticated, yet easier to use. Global computer networks are bringing state-of-the-art capabilities to all.



Current Opinion in Genetics and Development 1994, 4:383–388



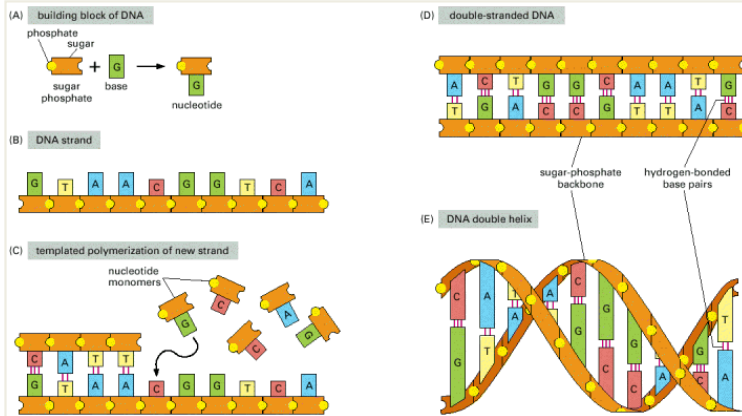
Table 1. Selected World-Wide Web (WWW) and Gopher servers.

Resource	WWW Uniform Resource Locator or Gopher address	Features and comments
World-Wide Web	http://info.cern.ch/	WWW project background; pointers to the world's on-line information; information on WWW software products; frequently asked questions (and answers).
NCSA Mosaic	http://www.ncsa.uiuc.edu/	Starting points for Internet exploration; Internet resources meta-index; "What's New" with NCSA and Internet.
Virtual Library: Biosciences	http://goji.harvard.edu/	Comprehensive list and links to WWW resources for biology and medicine.
ExPASy	http://expasy.hcuge.ch/	Comprehensive library of documents describing e-mail servers, databases and software for molecular biology.
NCBI GenBank®	http://ncbi.nlm.nih.gov/	Information for submitting and updating sequences; GenBank® release notes; homology and text searching of sequence databases; Entrez MEDLINE browser.
Genome Data Base (GDB)	http://gdbwww.gdb.org/	Search GDB and On-line Mendelian Inheritance in Man (OMIM)
CEPH/Généthon	http://www.genethon.fr/	Human genome genetic and physical mapping data; search for information on YAC clones and STS; QUICKMAP software for viewing CEPH map.
Cooperative Human Linkage Center (CHLC)	gopher://gopher.chlc.org/	Markers, genotype data and integrated maps.
UK Human Genome Mapping Project	gopher://menu.crc.ac.uk/	Primers, probes and chromosome abnormality database; cell lines from patients with genetic disorders or cytogenetic abnormalities.
MIT Genome Center	http://www-genome.wi.mit.edu/	Quarterly data releases from the human physical mapping and mouse genetic mapping projects.
Jackson Laboratory	http://www.jax.org/	Locus catalog and genetic maps of the mouse.
Dan's Favorite BioGopher	gopher://gopher.gdb.org/	Library catalogs around the world; world-wide campus phonebooks; search for people by name, location, research interest or funding agency.

GENOTYPE TO PHENOTYPE

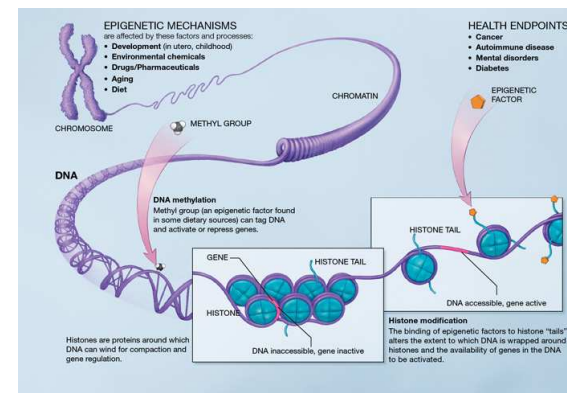
DNA

Alberts et al. 2002



Chromatin

NIH Common Fund <http://commonfund.nih.gov/epigenomics/figure.aspx>



Prediction of Chromatin States

Ernst et al. 2011

ARTICLE

doi:10.1038/nature09906

Mapping and analysis of chromatin state dynamics in nine human cell types

Jason Ernst^{1,2}, Pouya Kheradpour^{1,2}, Tarjet S. Mikkelsen¹, Noam Shores¹, Lucas D. Ward^{1,2}, Charles E. Epstein¹, Xiaolan Zhang¹, Li Wang¹, Robbyn Issner¹, Michael Coyne¹, Manolis Kellis^{1,3,4}, Timothy Durham¹, Manolis Kellis^{1,3,4} & Bradley E. Bernstein^{1,3,4}

Candidate state annotation	
Active promoter	
Weak promoter	
Inactive/poised promoter	
Strong enhancer	
Strong enhancer	
Weak/poised enhancer	
Weak/poised enhancer	
Insulator	
Transcriptional transition	
Transcriptional elongation	
Weak transcribed	
Polycomb repressed	
Heterochrom; low signal	
Repetitive/CNV	
Repetitive/CNV	

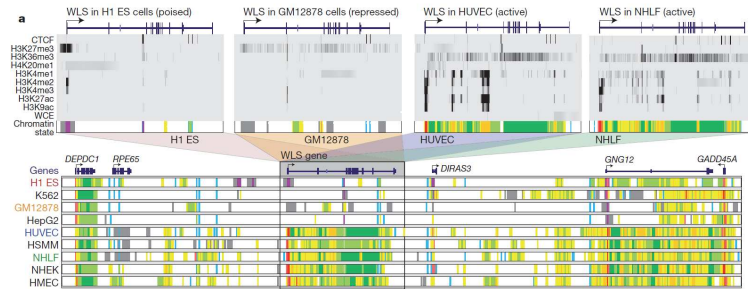


VANDERBILT UNIVERSITY



Prediction of Chromatin States

Ernst et al. 2011

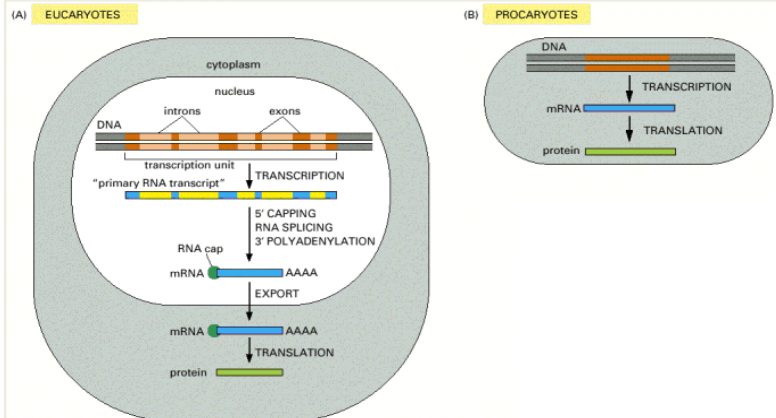


VANDERBILT UNIVERSITY



Transcription

Alberts et al. 2002

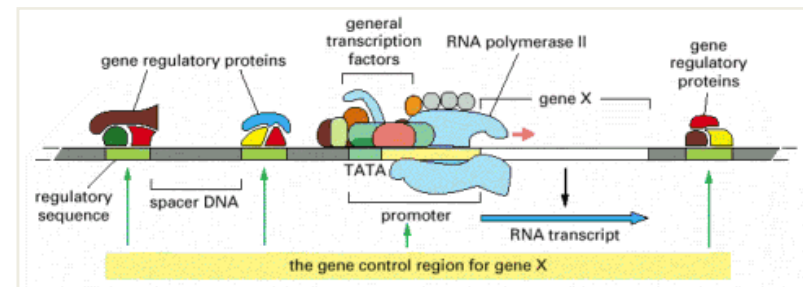


VANDERBILT UNIVERSITY



Regulation of Transcription

Alberts et al. 2002



VANDERBILT UNIVERSITY



Transcriptional Complexity

26 JUNE 2009 VOL 324 SCIENCE

Diversity and Complexity in DNA Recognition by Transcription Factors

Gwenael Badis,^{1*} Michael F. Berger,^{2,3*} Anthony A. Philippakis,^{2,3,4*} Shaheynoor Talukder,^{1,5*} Andrew R. Gehrke,^{2*} Savina A. Jaeger,^{2*} Esther T. Chan,^{5*} Genita Metzler,⁶ Anastasia Vedenko,⁷ Xiaoyu Chen,¹ Hanna Kuznetsov,⁶ Chi-Fong Wang,⁸ David Coburn,¹ Daniel E. Newburger,² Quaid Morris,^{1,5,9,10} Timothy R. Hughes,^{1,5,10}† Martha L. Bulyk,^{2,3,4,11}†



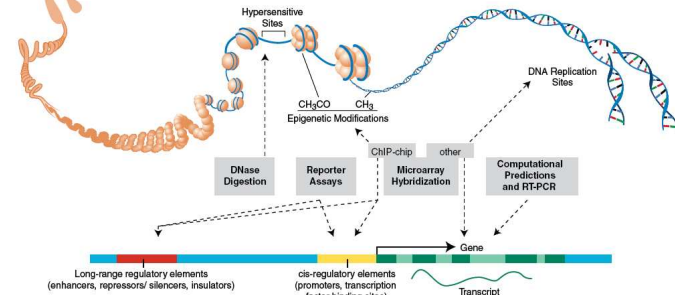
VANDERBILT UNIVERSITY



Transcriptional Complexity

The ENCODE (ENCyclopedia Of DNA Elements) Project

22 OCTOBER 2004 VOL 306 SCIENCE



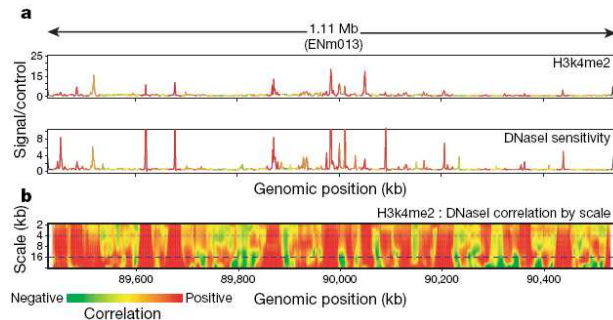
VANDERBILT UNIVERSITY



Identification and analysis of functional elements in 1% of the human genome by the ENCODE pilot project

The ENCODE Project Consortium*

NATURE | Vol 447 | 14 June 2007

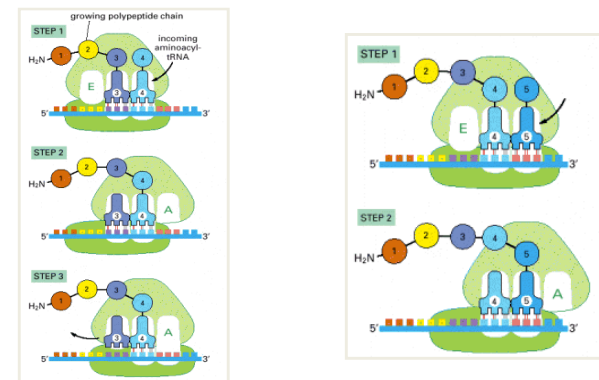


Negative Correlation Positive



Translation

Alberts et al. 2002

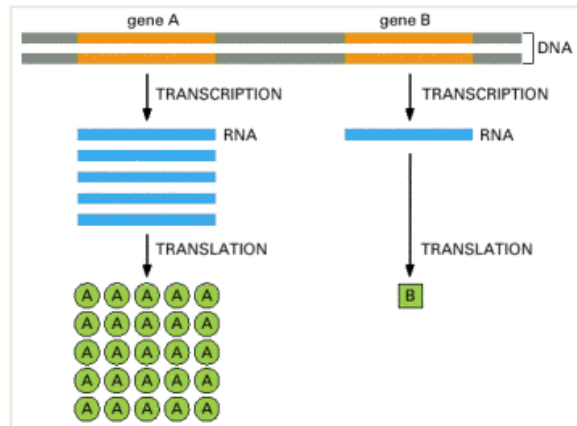


VANDERBILT UNIVERSITY



Translation

Alberts et al. 2002



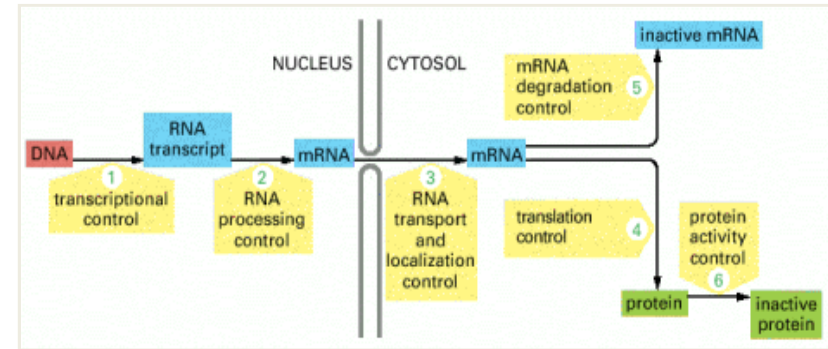
CHGR

VANDERBILT UNIVERSITY

WILLIAM S. BUSH PhD MS

Control of Gene Expression

Alberts et al. 2002



Alberts et al. 2002

CHGR

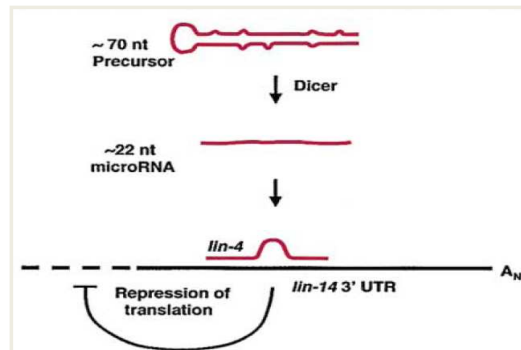
VANDERBILT UNIVERSITY

WILLIAM S. BUSH PhD MS

Cell, Vol. 107, 823-826, December 28, 2001, Copyright ©2001 by Cell Pr

microRNAs: Tiny Regulators with Great Potential

Victor Ambros¹
Department of Genetics
Dartmouth Medical School
Hanover, New Hampshire 03755



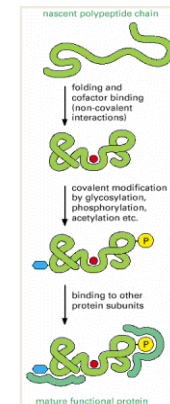
CHGR

VANDERBILT UNIVERSITY

WILLIAM S. BUSH PhD MS

Protein Folding

Alberts et al. 2012



CHGR

VANDERBILT UNIVERSITY

WILLIAM S. BUSH PhD MS

Amino Acids

Alberts et al. 2012

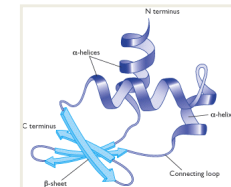
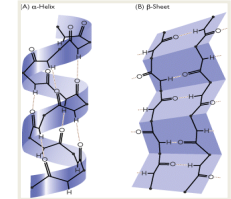
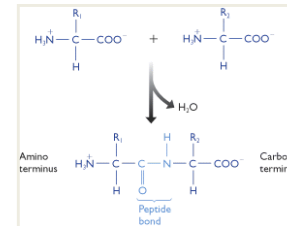
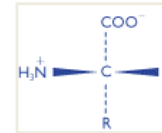
AMINO ACID	SIDE CHAIN	AMINO ACID	SIDE CHAIN
Aspartic acid Asp D	negative	Alanine Ala A	nonpolar
Glutamic acid Glu E	negative	Glycine Gly G	nonpolar
Arginine Arg R	positive	Valine Val V	nonpolar
Lysine Lys K	positive	Leucine Leu L	nonpolar
Histidine His H	positive	Isoleucine Ile I	nonpolar
Asparagine Asn N	uncharged polar	Proline Pro P	nonpolar
Glutamine Gln Q	uncharged polar	Phenylalanine Phe F	nonpolar
Serine Ser S	uncharged polar	Methionine Met M	nonpolar
Threonine Thr T	uncharged polar	Tryptophan Trp W	nonpolar
Tyrosine Tyr Y	uncharged polar	Cysteine Cys C	nonpolar

POLAR AMINO ACIDS

NONPOLAR AMINO ACIDS

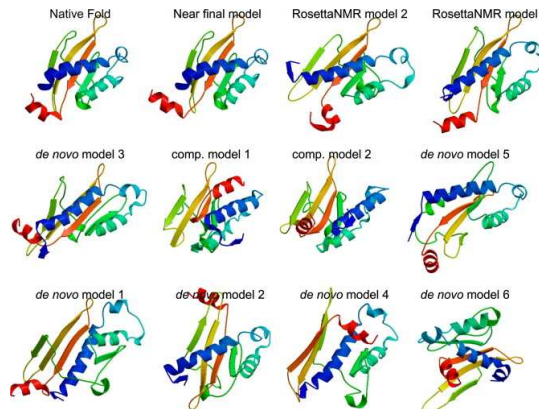
Amino Acids Peptides Proteins

Brown 2002



Protein Structural Prediction

Meiler and Baker 2005



indiana university Center for Computational Biology and Bioinformatics temple university Center for Information Science and Technology



DisProt News

Current release: 5.1
Release date: 05/28/2010
Number of proteins: 554
Number of disordered regions: 1223

Release notes

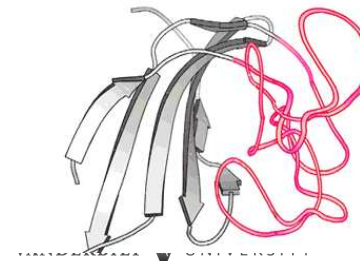
- Latest additions:
- Dynein intermediate chain, cytosolic [Isoform 5a (1)]
 - Fermitin family homolog 1
 - Translocase of chloroplast 159, chloroplastic
 - Inner membrane protein ALBINO3, chloroplastic [Isoform 1]
 - Dehydrin ERD14
 - more...

Download DisProt

Download DisProt in FASTA or XML format.

Database of Protein Disorder

The Database of Protein Disorder (DisProt) is a curated database that provides information about proteins that lack fixed 3D structure in their putatively native states, either in their entirety or in part. DisProt is a collaborative effort between Center for Computational Biology and Bioinformatics at Indiana University School of Medicine and Center for Information Science and Technology at Temple University.



www.disprot.org

Impact of Protein Changes

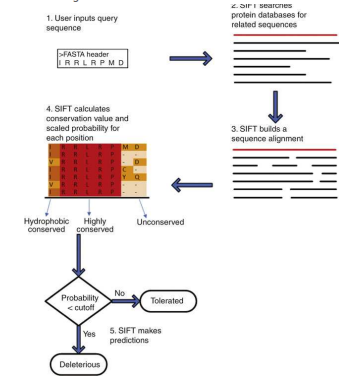
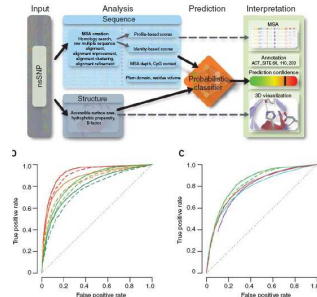
Adzhubei et al, Ng and Henikoff.

A method and server for predicting damaging missense mutations

Ivan A Adzhubei^{1,2}, Steffen Schmidt^{2,7}, Leonid Peshkin^{3,7}, Vasily E Ramensky⁴, Anna Gerasimova⁵, Peer Bork⁶, Alexey S Kondrashov⁵ & Shamil R Sunyaev¹

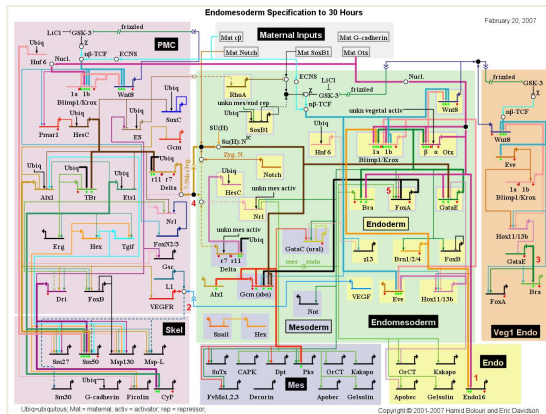
Predicting Deleterious Amino Acid Substitutions

Pauline C. Ng^{1,2} and Steven Henikoff^{1,3,4}

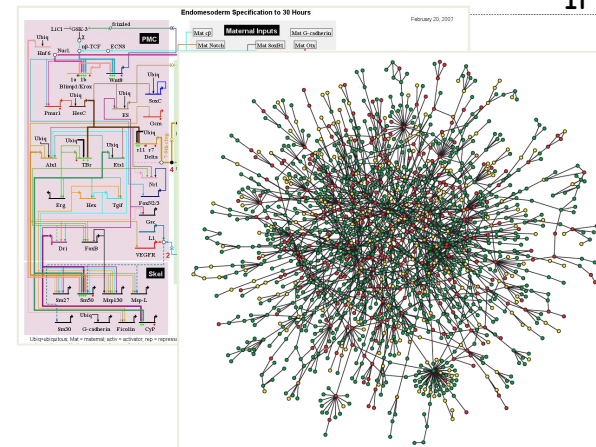


BIOMOLECULAR INTERACTIONS DRIVE BIOLOGY

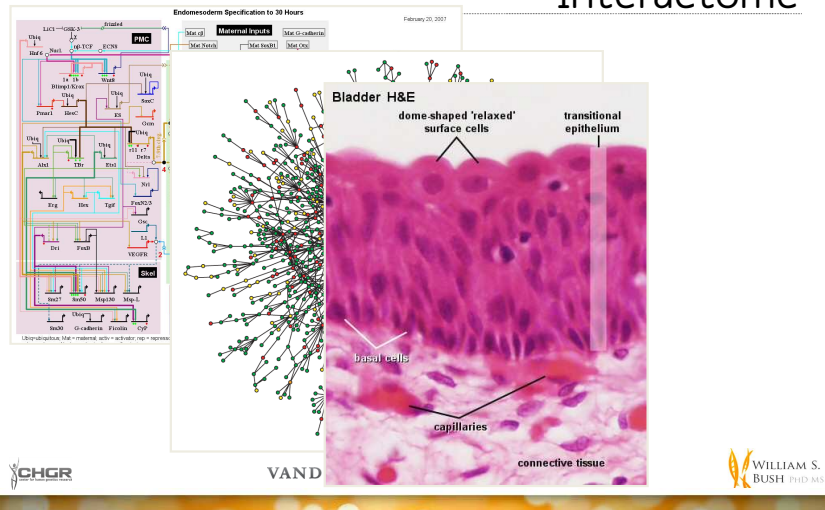
Interactome



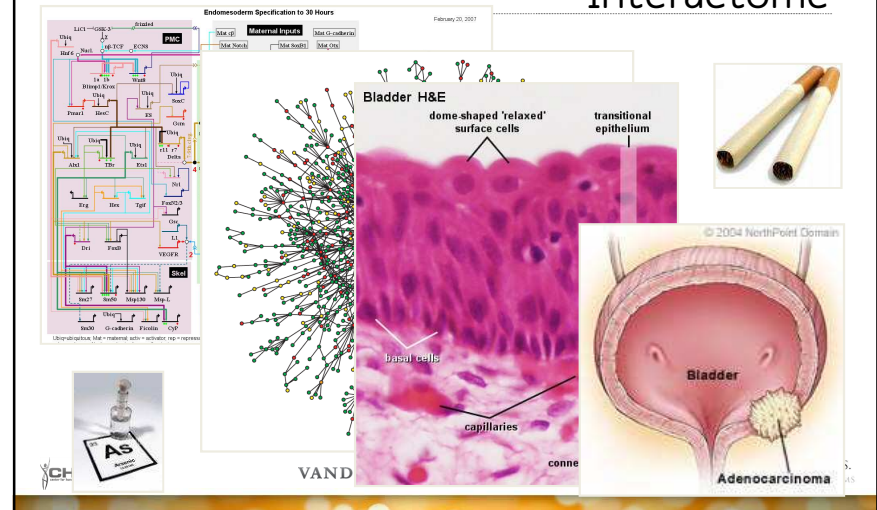
Interactome



Interactome



Interactome



SOME BIG QUESTIONS

Genetics and Life

THE PHENOGENETIC LOGIC OF LIFE

Kenneth M. Weiss

NATURE REVIEWS | GENETICS VOLUME 6 | JANUARY 2005 | 37

How does this...

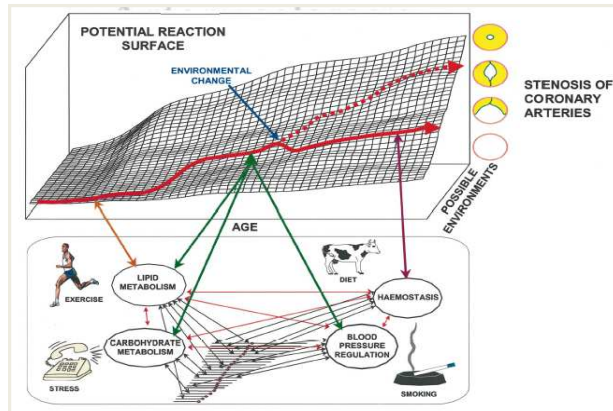
[illegible]

...become this



Genetic Architecture

Sing et al., *Arter. Thromb. Vasc. Biol.* (2003)



The Tree of Life



<http://tolweb.org/tree/>

Molecular Phylogenetics

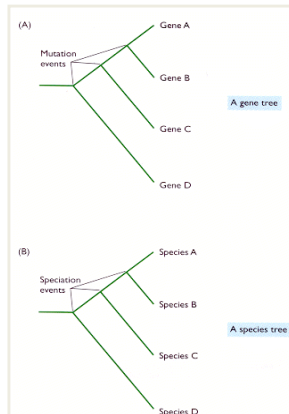
Multiple alignment

```

1 AGGCCAAGCCATAGCTGTCC
2 AGGCAAAGACATACCTGACC
3 AGGCCAAGACATAGCTGTCC
4 AGGCAAAGACATACCTGTCC
    
```

Distance matrix

	1	2	3	4
1	—	0.20	0.05	0.15
2		—	0.15	0.05
3			—	0.10
4				—



MEASURING DNA

EXCLUSIVE Q&A
DR. LAURA ON THE OFFENSIVE

CRACKING THE CODE!

THE INSIDE STORY OF HOW THESE BITTER RIVALS MAPPED OUR DNA

Q & A WITH DR. CRAIG VENTER AND FRANCIS COLLINS

Dr. Craig Venter

Francis Collins

The inside story of how these bitter rivals mapped our DNA. We have to find that elusive machine forever.

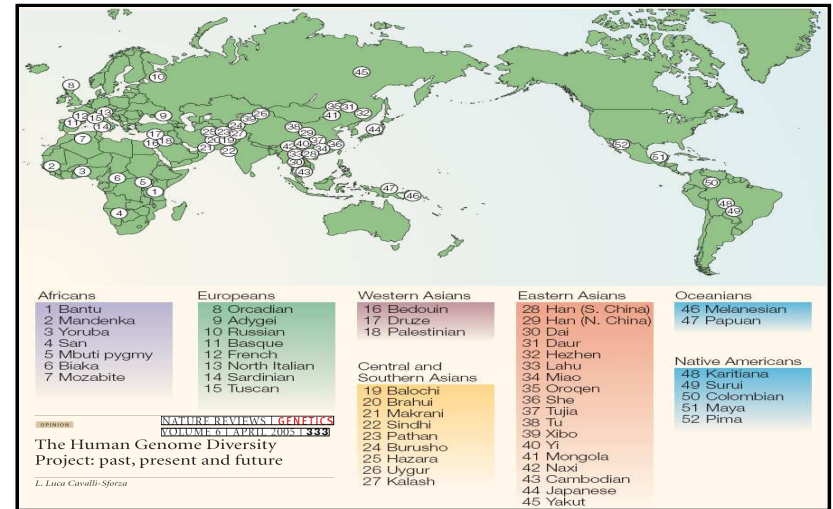
Cracking The Code!

CELEBRITY

VANDERBILT UNIVERSITY July 3, 2000



WILLIAM S.
BUSH PHD MS

VANDERBILT UNIVERSITY

38 WILLIAM S. BUSH PHD MS

Subject #1

Subject #2

Subject #3

-- AGCTCA --
-- AGCTCA --

Three *genotypes* (GG, GC, CC)

VANDERBILT UNIVERSITY

WILLIAM S.
BUSH PhD MS

- ~ 1 SNP every 100 bp
- ~ 30 million SNPs
- ~500,000 SNPs in coding DNA
 - Synonymous (silent)
 - Nonsynonymous
 - Deleterious effect
 - Beneficial effect
 - No effect

VANDERBILT UNIVERSITY

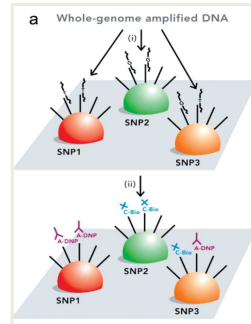
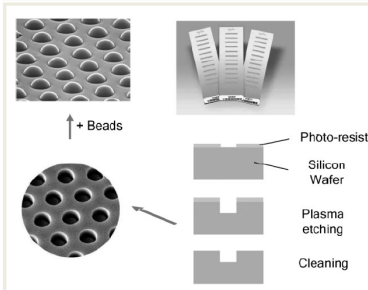
WILLIAM S.
BUSH PHD MS

Capturing SNPs

Whole genome genotyping technologies on the BeadArray™ platform

Frank J. Steemers and Kevin L. Gunderson
Illumina, Inc., San Diego, CA, USA

Biotechnol. J. 2007, 2, 41–49



BREAKTHROUGH OF THE YEAR

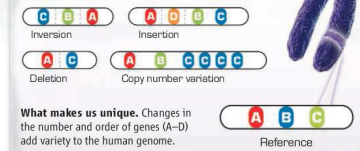
Science, December 21, 2007

Human Genetic Variation

Equipped with faster, cheaper technologies for sequencing DNA and assessing variation in genomes on scales ranging from one to millions of bases, researchers are finding out how truly different we are from one another

THE UNVEILING OF THE HUMAN GENOME ALMOST 7 YEARS AGO cast the first faint light on our complete genetic makeup. Since then, each new genome sequenced and each new individual studied has illuminated our genomic landscape in ever more detail. In 2007, researchers came to appreciate the extent to which our genomes differ from person to person and the implications of this variation for deciphering the genetics of complex diseases and personal traits.

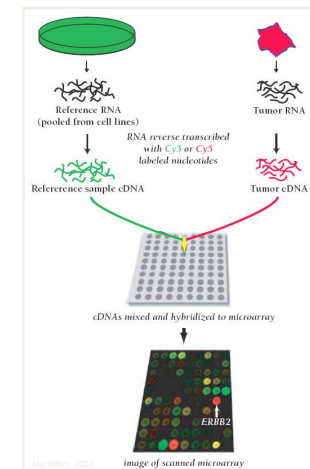
Less than a year ago, the big news was triangulating variation between us and our primate cousins to get a better handle on genetic changes along the evolutionary tree that led to humans. Now, we have moved from asking what in our DNA makes us human to striving to know what in my DNA makes me me.



What makes us unique. Changes in the number and order of genes (A–D) add variety to the human genome.

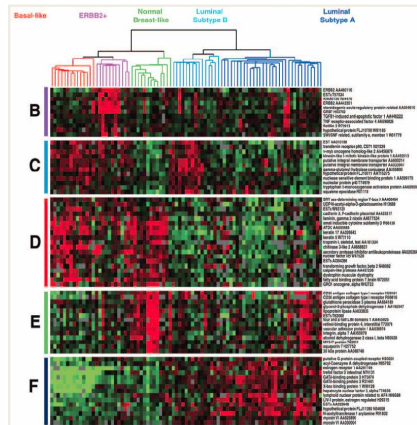
MEASURING RNA

cDNA Microarrays



Heatmaps

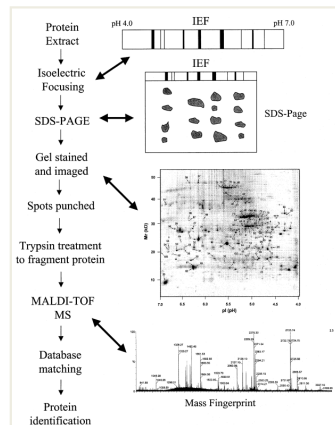
Mol Interv. 2002



MEASURING PROTEINS

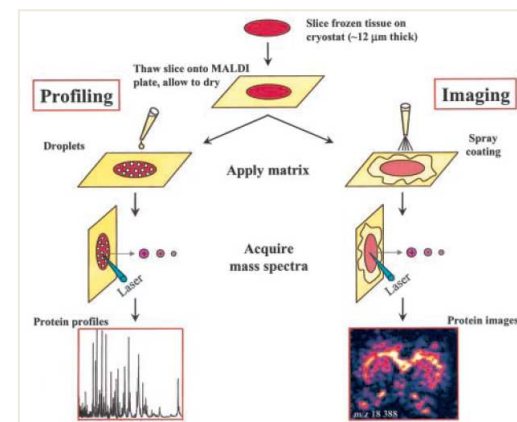
2D Gels and Mass Spectrometry

Metabolic Engineering 4, 98-106 (2002)

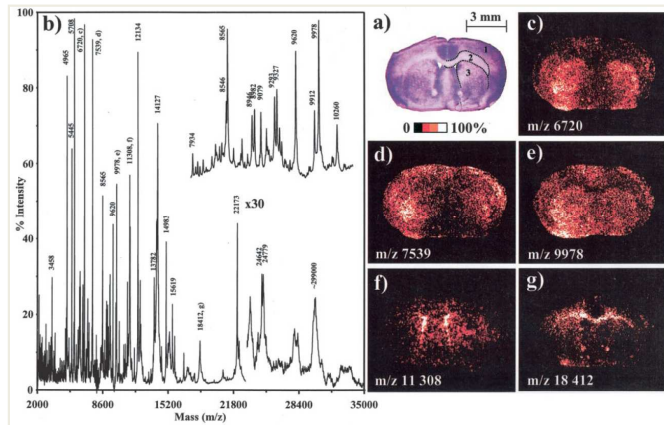


Protein Profiling in Tissues

Am J Pathol. 2004

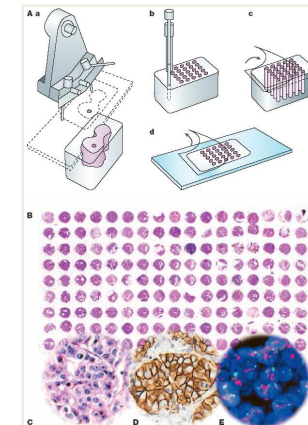


Protein Profiling in Tissues



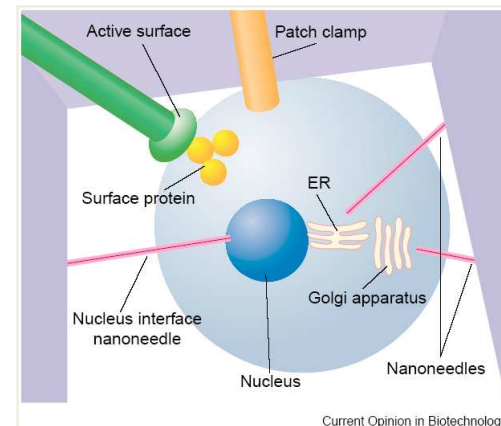
Tissue Microarrays

Nat Rev Drug Discov. 2003



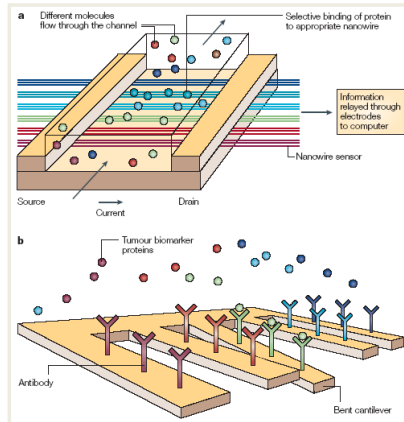
EMERGING TECHNOLOGIES

Nanotechnology



Nanotechnology

Nature Reviews Cancer 2005

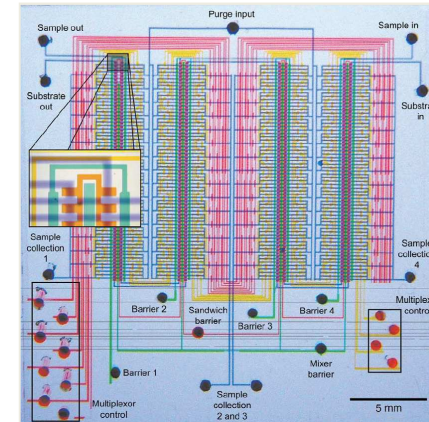


VANDERBILT UNIVERSITY



Lab-on-a-Chip

Current Opinion in Structural Biology 2003



VANDERBILT UNIVERSITY



The \$1000 Genome

J. Craig Venter

INSTITUTE



Press Release

FOR IMMEDIATE RELEASE

J. CRAIG VENTER SCIENCE FOUNDATION ANNOUNCES \$500,000 TECHNOLOGY PRIZE FOR ADVANCES LEADING TO THE \$1,000 HUMAN GENOME

ROCKVILLE, MD (September 23, 2003). The J. Craig Venter Science Foundation announced today a \$500,000 Genomic Technology Prize. The prize, to be awarded one time only, is aimed at stimulating the scientific and technology research community to significantly advance automated DNA sequencing so that a human genome can be sequenced for \$1,000 or less as soon as possible. The prize was announced during New Frontiers in Sequencing Technology session at the 15th annual Genome Sequencing and Analysis Conference (GSAC) in Savannah, Georgia.

100 Genomes in 10 Days

ARCHON GENOMICS XPRIZE

DONATE | X PRIZE FOUNDATION

ARCHON X PRIZE FOR GENOMICS | TEAMS | NEWS & EVENTS | TAKE ACTION | DISCOVER | ABOUT

Revolution Through Competition.

TAKE ACTION

ARCHON X PRIZE FOR GENOMICS

- Introduction
- Why Genomics?
- Prize Overview
- The Promise of Personalized Medicine
- Frequently Asked Questions

PRIZE OVERVIEW

A \$10 MILLION PRIZE
FOR THE FIRST TEAM TO SUCCESSFULLY SEQUENCE
100 HUMAN GENOMES IN 10 DAYS

What Inspired this X PRIZE?

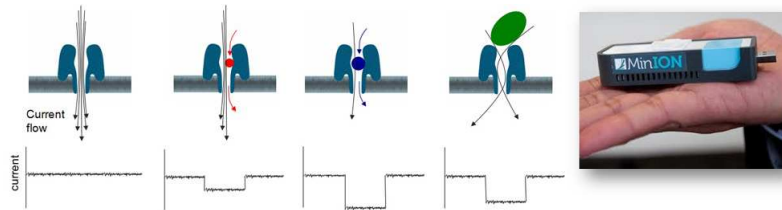
In 2000, Dr. J. Craig Venter led the first private team to successfully sequence a complete human genome. In the preceding decades combined governmental and private funding efforts spent \$100s of millions to develop the instrumentation required. It took the Venter team \$100 million and nine months to achieve their historic accomplishment.

"The X PRIZE Foundation is taking a bold and important step — incentivizing innovation in technology that will alleviate a bottleneck in genomic sciences. We, people who suffer from genetic conditions, will benefit so much from

The J. Craig Venter Science Foundation offered the \$500,000 Innovation in Genomics Science and Technology Prize in September 2003 aimed at stimulating development of less expensive and faster sequencing technology. To attract even more resources to this exceptionally worthy goal, Dr. Venter joined forces with the X PRIZE Foundation, wrapping his competition and prize purse into the Archon X PRIZE for Genomics.

Nanopore Technology

<http://www.nanoporetech.com>



VANDERBILT UNIVERSITY



1000 Genomes Project

1000 Genomes

A Deep Catalog of Human Genetic Variation

Home About Data Analysis Participants Contact Browser Wiki

LATEST ANNOUNCEMENTS

March 2010 Data Release

31 MARCH 2010
Final release of pilot project SNP calls

The final set of SNPs from Pilots 1, 2 and 3 are now available in VCF format. All 1000 Genomes pilot project files reference the NCBI build 36 assembly of the human genome.

Data access links: [EBI](#) / [NCBI](#)

Link to additional information: [README file](#)



VANDERBILT UNIVERSITY

www.1000genomes.org



Our other genomes

<http://nihroadmap.nih.gov/hmp/>

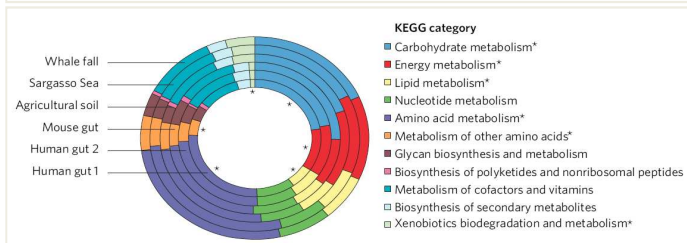
INSIGHT FEATURE

NATURE | Vol 449 | 18 October 2007 | doi:10.1038/nature06244

The Human Microbiome Project

Peter J. Turnbaugh, Ruth E. Ley, Micah Hamady, Claire M. Fraser-Liggett, Rob Knight & Jeffrey I. Gordon

A strategy to understand the microbial components of the human genetic and metabolic landscape and how they contribute to normal physiology and predisposition to disease.



VANDERBILT UNIVERSITY



DATABASES

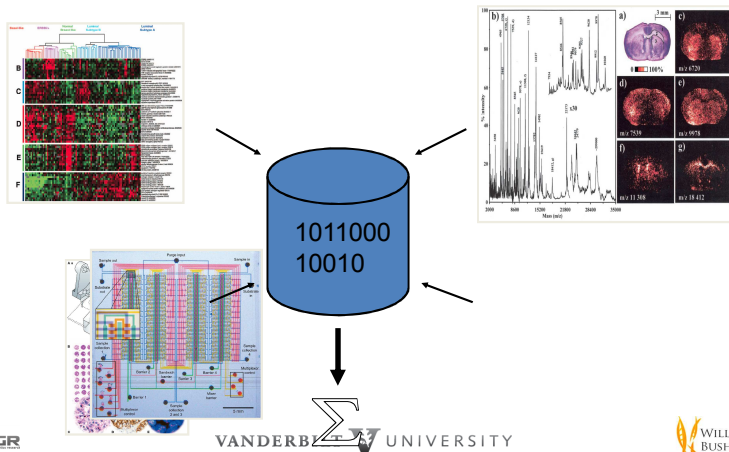


VANDERBILT UNIVERSITY

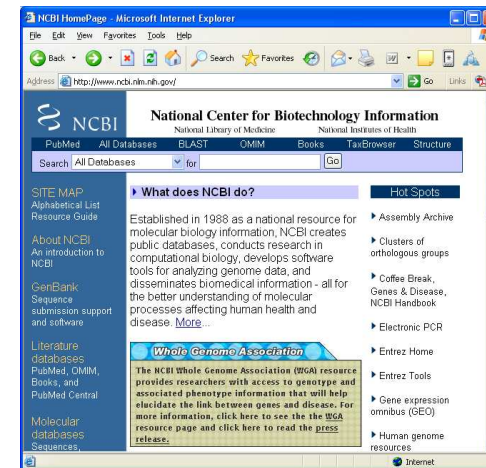
60



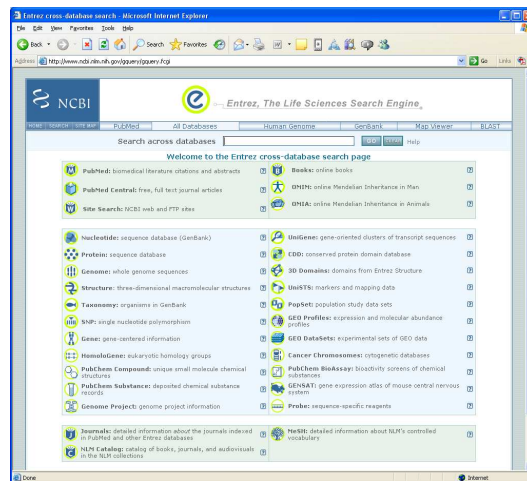
Bioinformatics: Databases



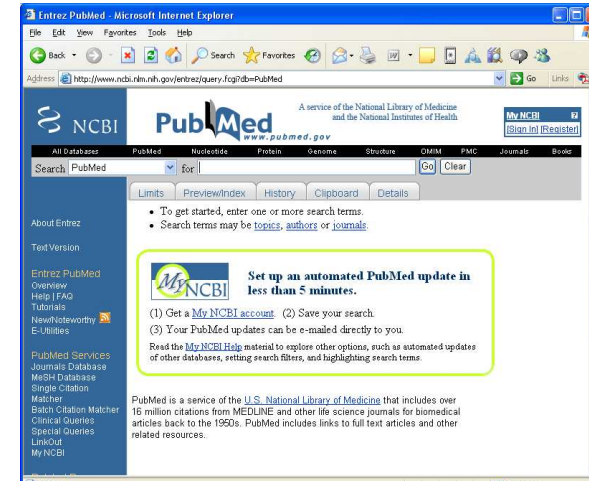
<http://www.ncbi.nlm.nih.gov>



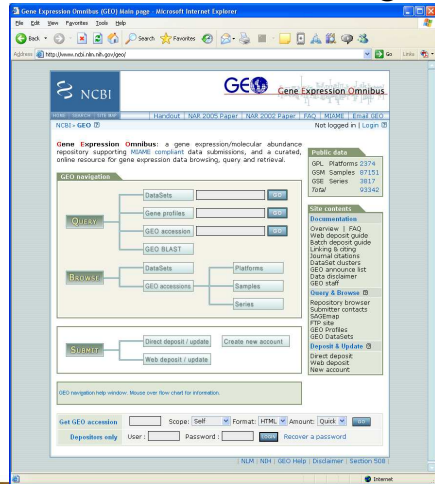
<http://www.ncbi.nlm.nih.gov/gquery/gquery.fcgi>



<http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?db=PubMed>

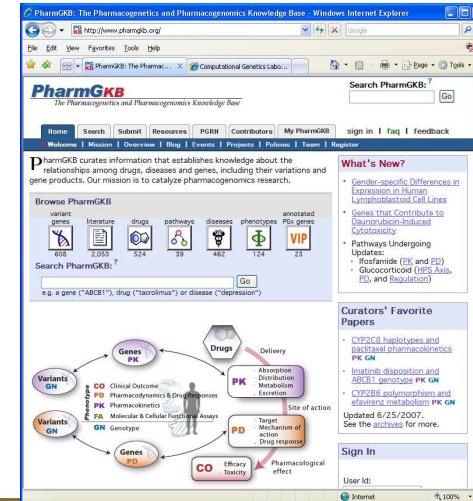


<http://www.ncbi.nlm.nih.gov/geo/>



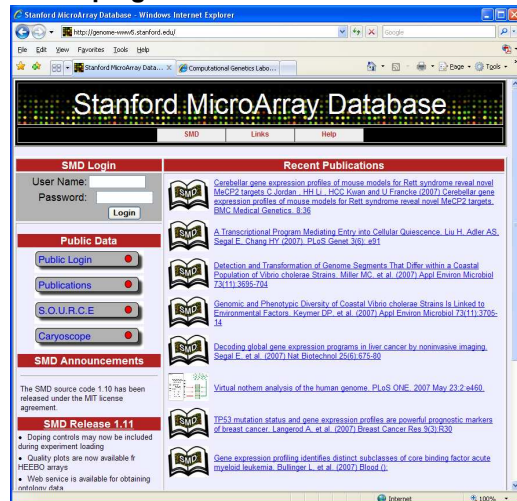
WILLIAM S. BUSH PhD MS

www.pharmgkb.org



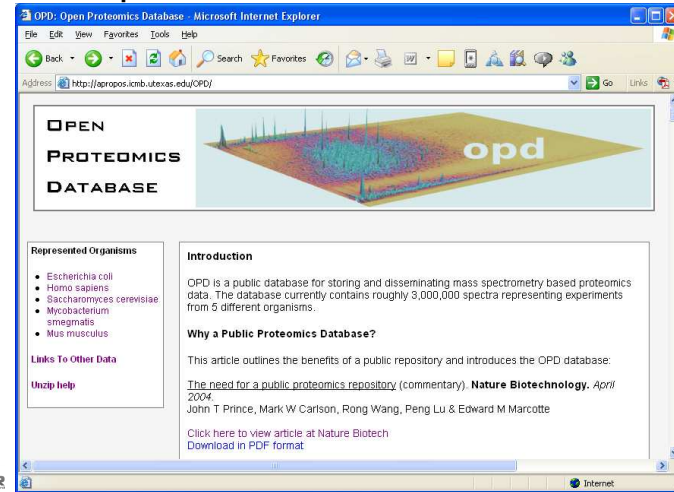
WILLIAM S. BUSH PhD MS

<http://genome-www5.stanford.edu/>



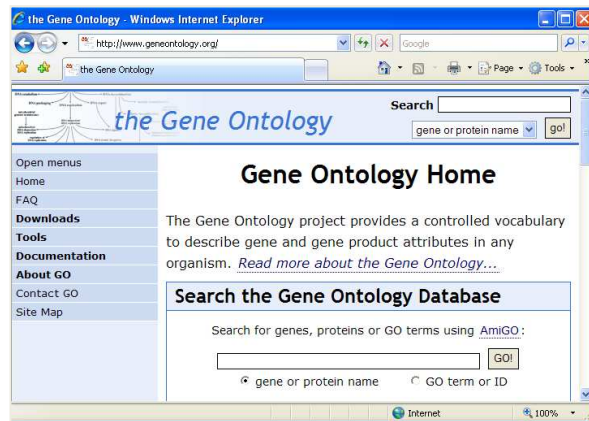
WILLIAM S. BUSH PhD MS

<http://bioinformatics.icmb.utexas.edu/OPD/>



WILLIAM S. BUSH PhD MS

<http://www.geneontology.org>



VANDERBILT UNIVERSITY



<http://www.geneontology.org>

```

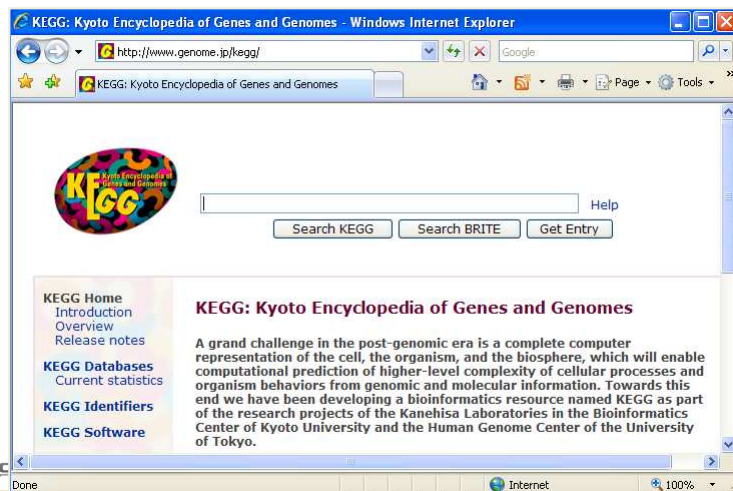
+ all : all [218539]
+ O:0005575 : cellular_component [152301]
+ O:0005623 : cell [107815]
+ O:0044464 : cell part [107776]
+ O:0005622 : intracellular [86035]
+ O:0044424 : intracellular part [85270]
+ O:0005737 : cytoplasm [72113]
+ O:0044444 : cytoplasmic part [66867]
+ O:0005794 : Golgi apparatus [1618]
    
```



VANDERBILT UNIVERSITY

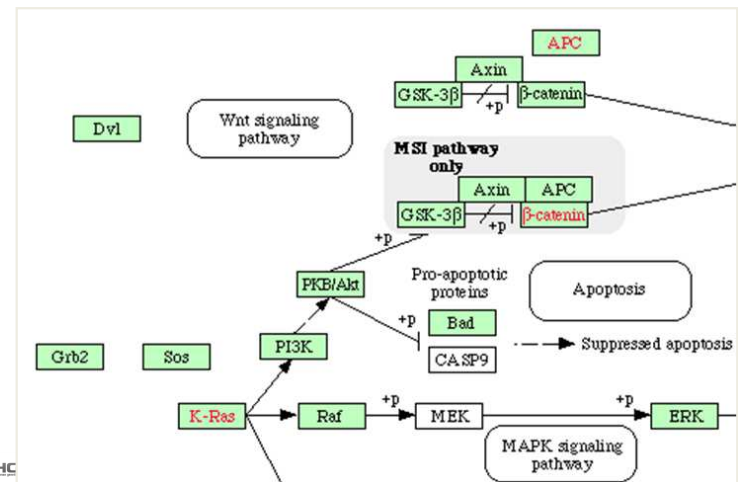


<http://www.genome.jp/kegg/>



AM S.
PHD MS

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



LIAM S.
PHD MS

<http://string-db.org/>

Home • Download • Help/Info

STRING 8.3

STRING - Known and Predicted Protein-Protein Interactions

search by name | search by protein sequence | multiple names | multiple sequences

protein name: (examples: #1 #2 #3)

(STRING understands a variety of protein names and accessions; you can also try a [random entry](#))

organism: auto-detect

interactors wanted: COGs Proteins

Reset GO!

please enter your protein of interest...

What it does ...

STRING is a database of known and predicted protein interactions. The interactions include direct (physical) and indirect (functional) associations; they are derived from four sources:

Genomic Context | High-throughput Experiments | (Conserved) Coexpression | Previous Knowledge

STRING quantitatively integrates interaction data from these sources for a large number of organisms, and transfers information between these organisms where applicable. The database currently covers 2,590,259 proteins from 630 organisms.

More Info | Funding / Support | Acknowledgements | Use Scenarios

STRING (Search Tool for the Retrieval of Interacting Genes/Proteins) is being developed at CPB, EMBL, SIB, KU, TUD and UZH. STRING references: Jensen et al. 2009 / 2007 / 2005 / 2003 / Snel et al. 2000. Miscellaneous: [Access Statistics](#), [Robot Access Guide](#), [STRING/STITCH Blog](#), [Supported Browsers](#).

What's New? This is version 8.3 of STRING - June 2010; the latest interaction data, updated textmining, and bugfixes ... **Sister Projects:** check out [STITCH](#) and [eggNOG](#) - two sister projects built on STRING data! **Previous Releases:** Trying to reproduce an earlier finding? Confused? Refer to our [old releases](#).

CHGR VANDERBILT UNIVERSITY WILLIAM S. BUSH PHD MS

ANALYSIS

CHGR VANDERBILT UNIVERSITY WILLIAM S. BUSH PHD MS

Mining Biomolecular Patterns

- Can we classify and/or predict biological and clinical endpoints using genetic, genomic, and/or proteomic data?
- Which biomolecules are important?
- What is their pattern or statistical relationship?

CHGR VANDERBILT UNIVERSITY WILLIAM S. BUSH PHD MS

Objectives

Data → Variable Selection → Statistical Modeling → Prediction Classification

$I = a_1x_1 + a_2x_2$

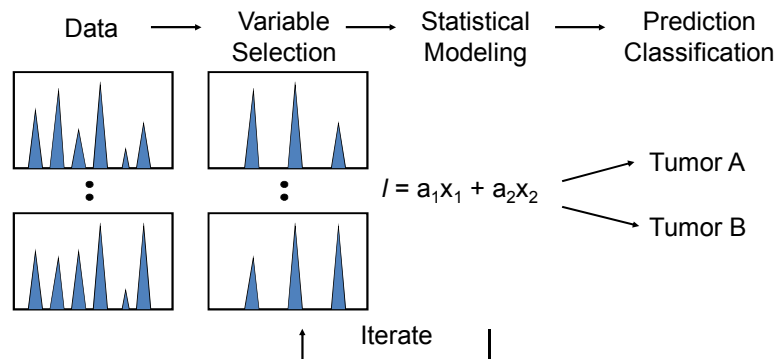
Tumor A

Tumor B

Iterate

CHGR VANDERBILT UNIVERSITY WILLIAM S. BUSH PHD MS

Objectives



Hypothesis Testing

"The truth is out there"

		Truth		
		H_0 False	H_0 True	
Decision	Reject H_0	Yes! Type I Error	Type I Error	Type I error Type II error Type III error
	Accept H_0	Type II Error	Yes!	

Permutation Testing

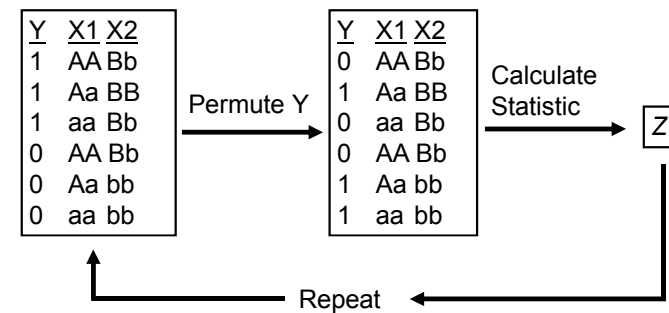
P. Good, Permutation Tests: A Practical Guide to Resampling Methods for Testing Hypotheses (2000)

- Many data-driven methods are nonparametric and model-free.
- Permutation testing can be used to assess statistical significance to allow formal hypothesis testing.
- Basic Idea: Randomize data so it is consistent with null hypothesis.

Permutation Testing

P. Good, Permutation Tests: A Practical Guide to Resampling Methods for Testing Hypotheses (2000)

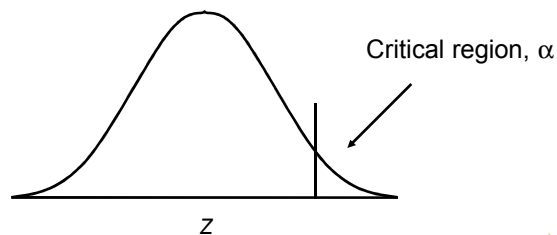
Example



Permutation Testing

P. Good, *Permutation Tests: A Practical Guide to Resampling Methods for Testing Hypotheses* (2000)

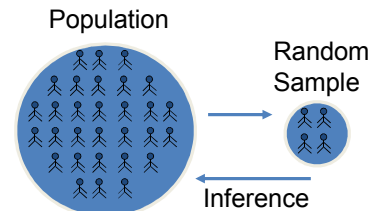
Distribution of Statistic under the Null Hypothesis from Many Permutations



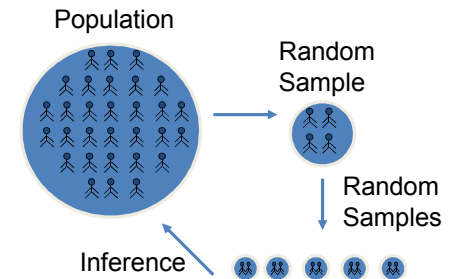
Bootstrapping

B. Efron, *Annals of Statistics* 7:1-26 (1979)
AC Davidson and DV Hinkley, *Bootstrap Methods and their Application* (1997)
CE Lunneborg, *Data Analysis by Resampling: Concepts and Applications* (2000)

Distribution Known



Distribution Unknown



Health is a Complex System

Sing et al., *Arter. Thromb. Vasc. Biol.* (2003)

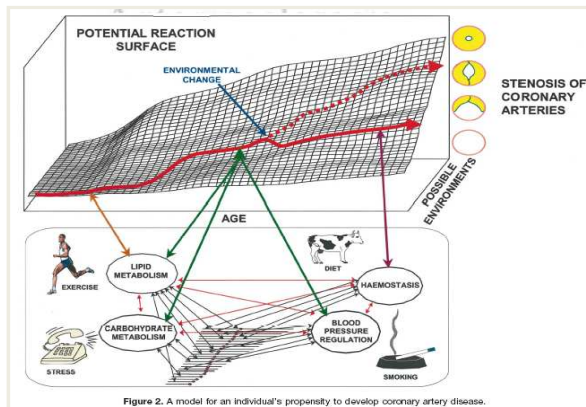


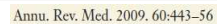
Figure 2. A model for an individual's propensity to develop coronary artery disease.

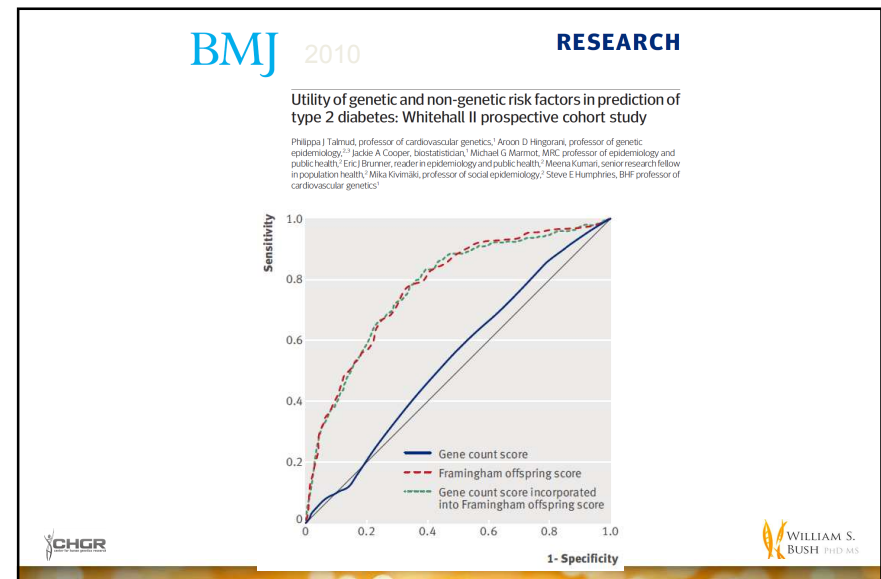
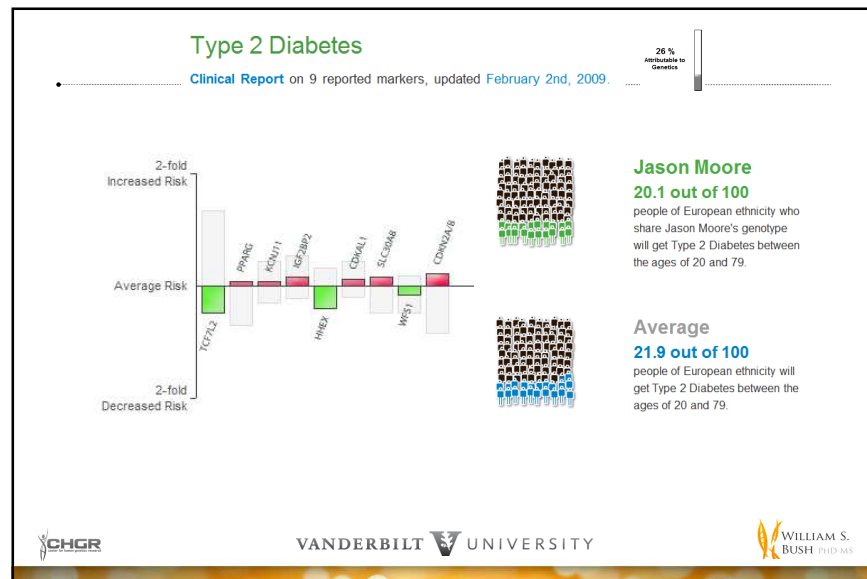
ARTICLES

Genome-wide association study of 14,000 cases of seven common diseases and 3,000 shared controls

The Wellcome Trust Case Control Consortium*

- 50 Research Groups
- 14,000 cases and 3,000 shared controls
- 500,000 SNPs
- Seven complex human diseases:
 - bipolar disorder (BD)
 - coronary artery disease (CAD)
 - Crohn's disease (CD)
 - hypertension (HT)
 - rheumatoid arthritis (RA)
 - type 1 diabetes (T1D)
 - type 2 diabetes (T2D)

LIAM S.
PHD MSVANDERBILT  UNIVERSITY



A Genetics Company Fails, Its Research Too Complex

By NICHOLAS WADE
Published: November 17, 2009

The New York Times

DeCode Genetics, a pioneering company that used the Icelandic population as its guinea pigs in detecting disease-causing mutations, filed for bankruptcy on Tuesday.

[Enlarge This Image](#)

A lab at deCODE Genetics, which found that the genetic nature of human disease was far more complex than anyone thought.

The company's demise suggests that the medical promise of the human genome may take much longer to be fulfilled than its sponsors had hoped. Based in Reykjavik, Iceland, it was founded in 1996 by Dr. Kari Stefansson, a research neurologist who worked at the University of Chicago and at Harvard. After the human genome sequence was achieved in 2003, Dr. Stefansson quickly realized that Iceland's excellent medical records, combined with the genealogical information available on its close-knit population, provided a fine test bed for seeking the roots of genetically complex diseases like cancer, diabetes and schizophrenia.

Related
Times Topics: deCODE Genetics, Incorporated

CHGR VANDERBILT UNIVERSITY WILLIAM S. BUSH PhD MS

VIEWPOINT

NATURE REVIEWS | GENETICS
VOLUME 11 | JUNE 2010

Missing heritability and strategies for finding the underlying causes of complex disease

Evan E. Eichler, Jonathan Flint, Greg Gibson, Augustine Kong, Suzanne M. Leal, Jason H. Moore and Joseph H. Nadeau

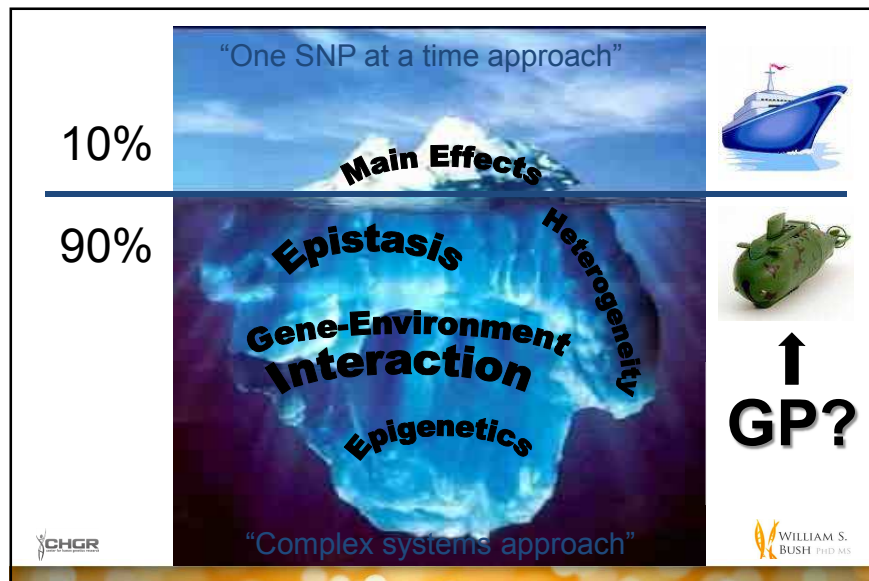
Q How should we solve the problem of 'missing heritability' in complex diseases?

Jason H. Moore. The case of the missing heritability for common human diseases should not be a mystery to anyone given the inherent complexity of the relationship between genotype and phenotype.

The time is now to philosophically and analytically retool for a complex genetic architecture or we will continue to underdeliver on the promises of human genetics.

Indeed, life, and thus genetics, is complicated⁴⁶ and some will soon ask, as seismologists have⁴⁷, whether we are trying to predict the unpredictable.

CHGR VANDERBILT UNIVERSITY WILLIAM S. BUSH PhD MS



Challenges

- Complexity – nonlinearity, heterogeneity
- Dimensionality – multiple genetic risk factors
- Scale – millions of attributes



CHGR

VANDERBILT UNIVERSITY

94

WILLIAM S. BUSH PhD MS



Using Expert Knowledge in GP

- Multi-Objective Fitness (GPTP'06)
 - Fitness = $F(\text{Accuracy} + \text{Knowledge})$
- Recombination (EvoBIO'07)
 - Recombine trees with good building blocks
- Mutation (PRIB'07)
 - Mutate trees with poor building blocks
- Sensible Initialization (CEC'09)
 - Initialize trees with good building blocks

CHGR

VANDERBILT UNIVERSITY

WILLIAM S. BUSH PhD MS

GPTP'07

Does Complexity Matter?

GUIDELINES

From artificial evolution to computational evolution:
a research agenda

Wolfgang Banzhaf, Guillaume Beslon, Steffen Christensen, James A. Foster,
François Képès, Virginie Lefort, Julian F. Miller, Miroslav Radman and
Jeremy J. Ramsden

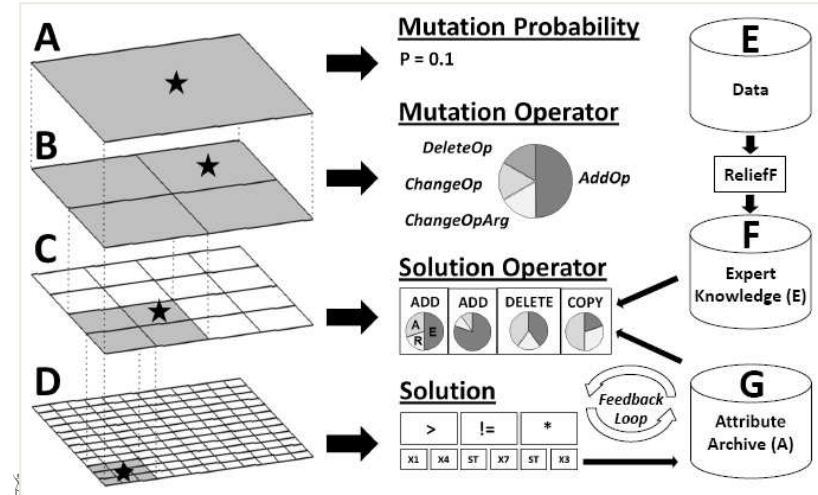
NATURE REVIEWS | **GENETICS** VOLUME 7 | SEPTEMBER 2006



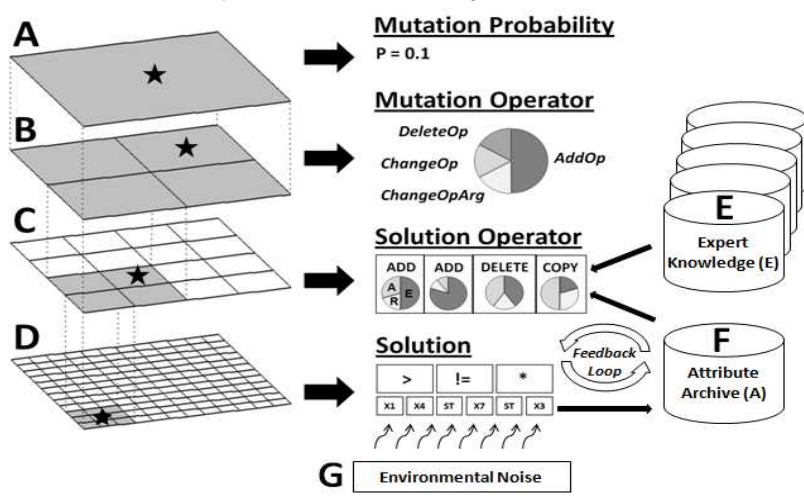
VANDERBILT UNIVERSITY



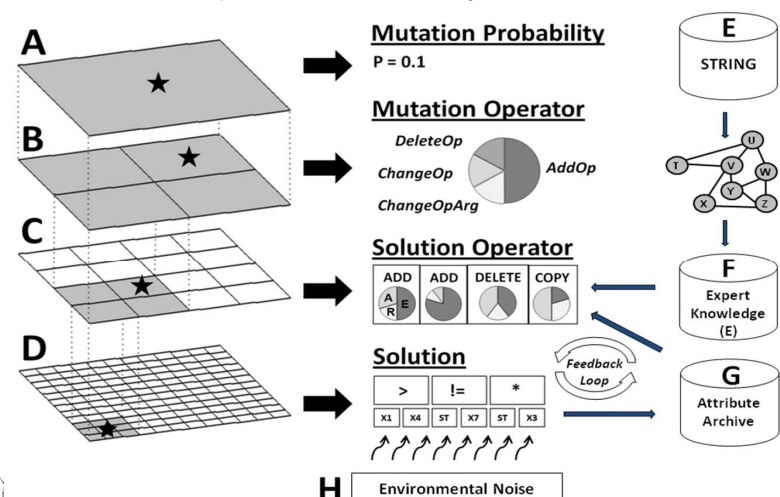
A Computational Evolution System – GPTP'08



A Computational Evolution System – GPTP'09



A Computational Evolution System – GPTP'10



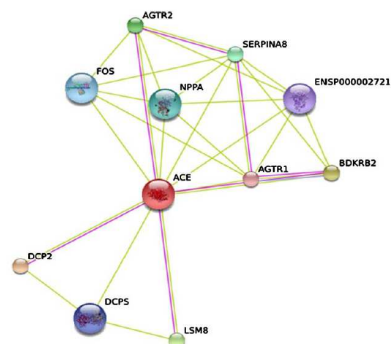
REVIEW

Exploiting the proteome to improve the genome-wide genetic analysis of epistasis in common human diseases

Kristine A. Pattin · Jason H. Moore

STRING

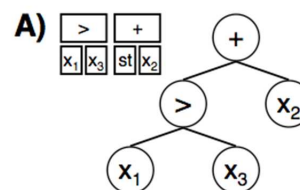
- 2.5 million proteins
- 630 organisms



CHGR

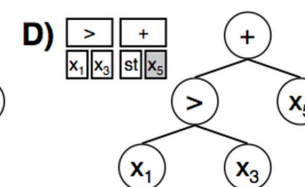
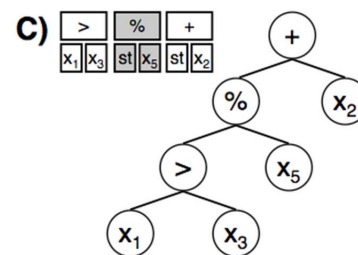
VANDER

WILLIAM S. BUSH PHD MS



B)

	x_1	x_2	x_3	x_4	x_5
x_1	0	.71	.84	.10	0
x_2	.63	0	.68	0	0
x_3	.84	.71	0	.10	.97
x_4	.24	0	.10	0	.32
x_5	0	0	.81	.32	0



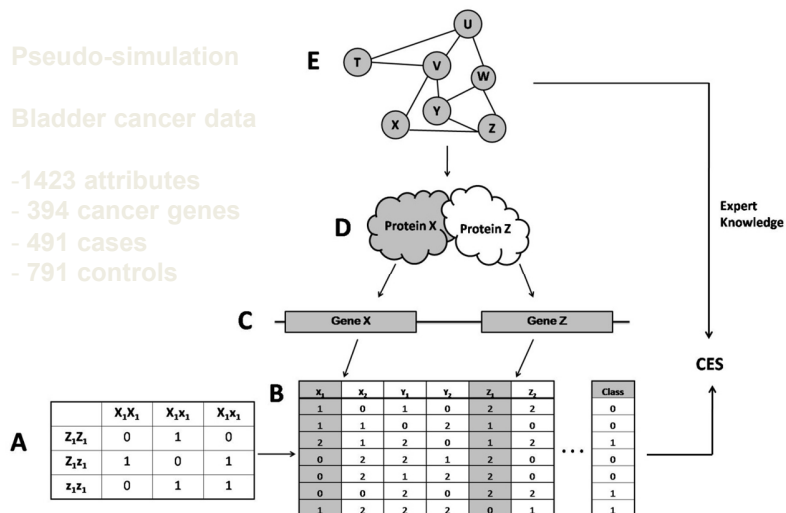
CHI

WILLIAM S. BUSH PHD MS

Pseudo-simulation

Bladder cancer data

- 1423 attributes
- 394 cancer genes
- 491 cases
- 791 controls



Experimental Design and Analysis

Datasets = 100 per model

Strong, medium and weak P x P interaction

Runs = 100 per dataset

Generations = 1000

Grid Size = 18x18

Control 1: no expert knowledge

Control 2: attributes with weak P x P interaction

Measured success in finding correct attributes

CHGR

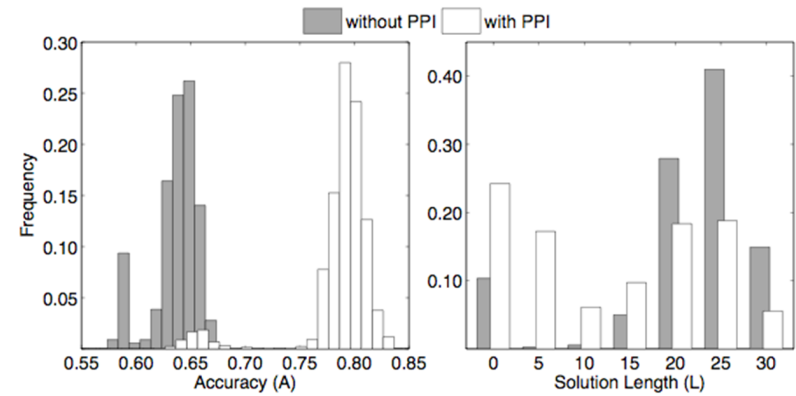
VANDERBILT UNIVERSITY

WILLIAM S. BUSH PHD MS

Results

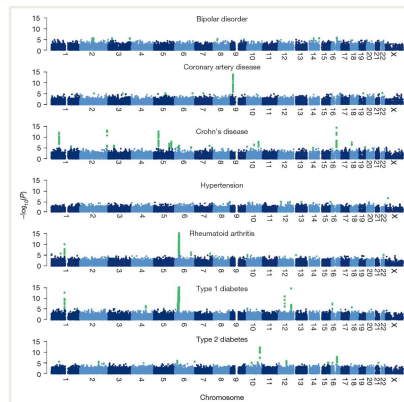
Table 1-1. Percentage of datasets in which CES successfully identified the correct SNP pairings as the most frequent, for the four confidence score scenarios considered.

Confidence Score of PPI	With Expert Knowledge	Control
0.998	100%	0%
0.963	100%	0%
0.933	96%	0%
0.916	80%	0%



Summary

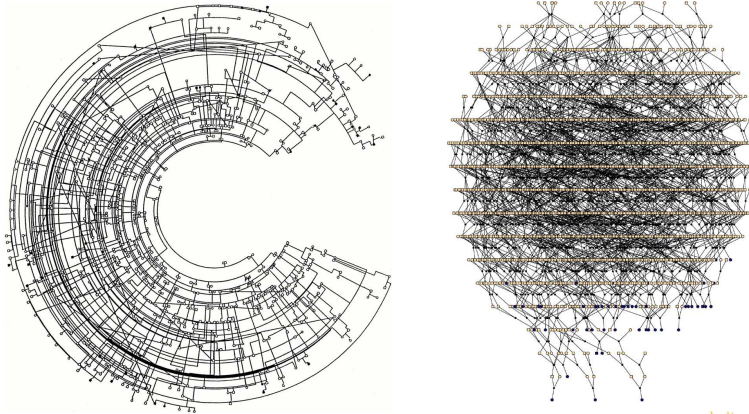
- Our CES is able to learn to exploit protein-protein interactions as a source of domain-specific knowledge.
- Can this scale to 10^6 or more attributes?



PROBLEM REPRESENTATION

Graph Partitioning

Racette et al. Neurology 2002, Lie et al., 2008



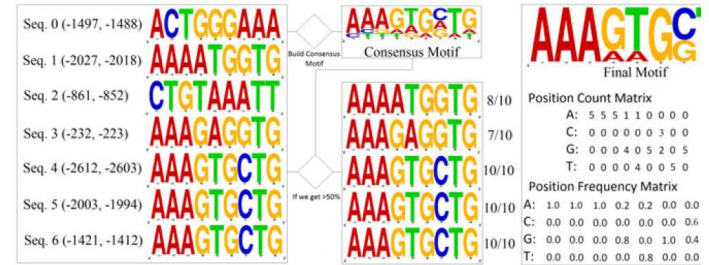
CHGR

VANDERBILT UNIVERSITY

WILLIAM S. BUSH PhD MS

Motif Discovery

Gonzalez-Alvarez et al. EvoBIO 2012



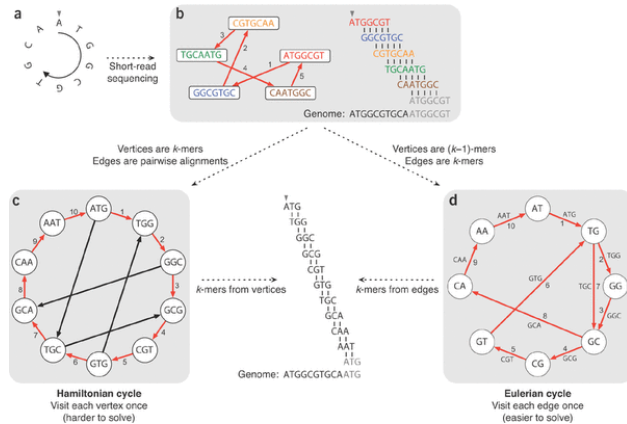
CHGR

VANDERBILT UNIVERSITY

WILLIAM S. BUSH PhD MS

Sequence Alignment

Compeau et al. Nature Biotech 2011



CHGR

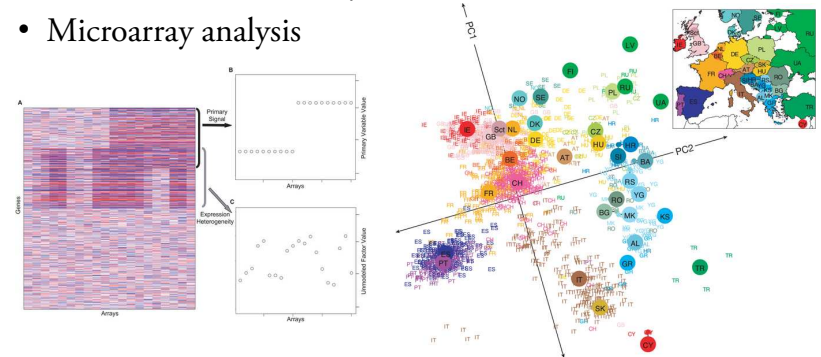
VANDERBILT UNIVERSITY

WILLIAM S. BUSH PhD MS

Latent Variable Discovery

Leek and Storey, PLoS Genetics 2007, Novembre et al, Nature 2008

- PCA for ethnic ancestry
- Microarray analysis



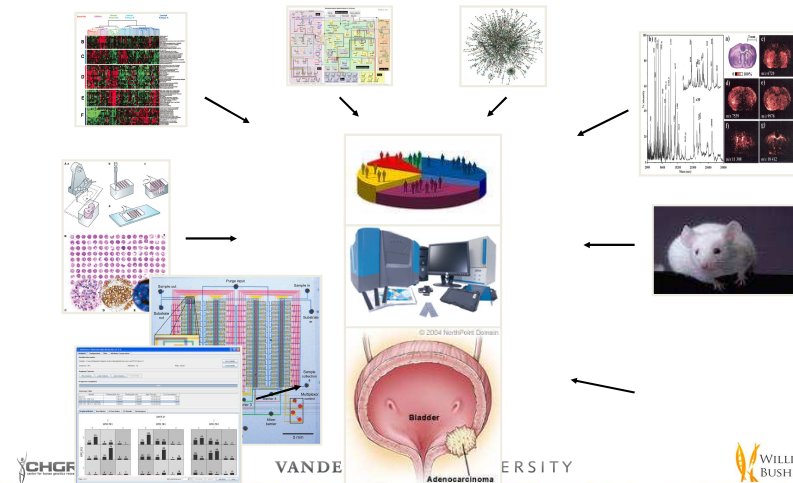
CHGR

VANDERBILT UNIVERSITY

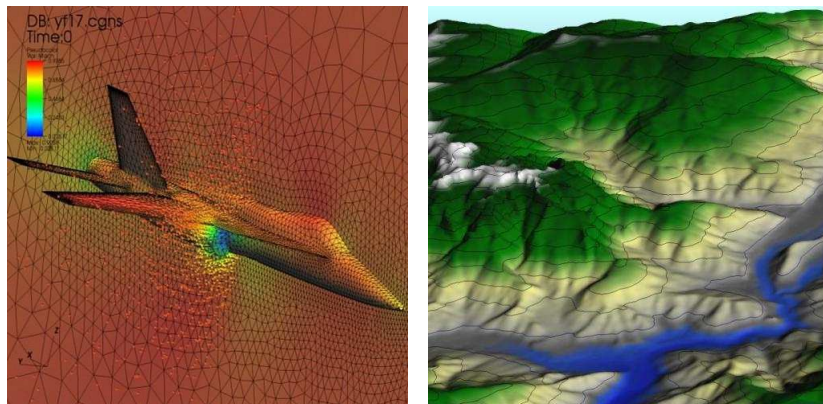
WILLIAM S. BUSH PhD MS

THE FUTURE

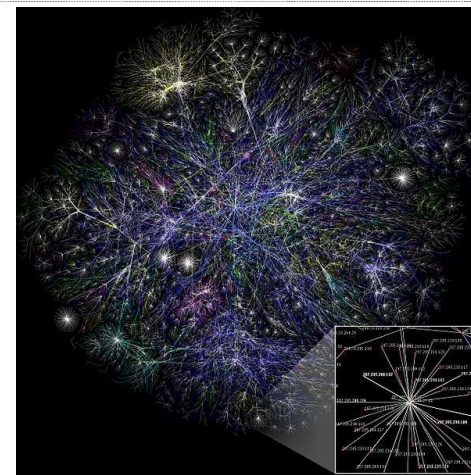
Collaborative Approach Needed



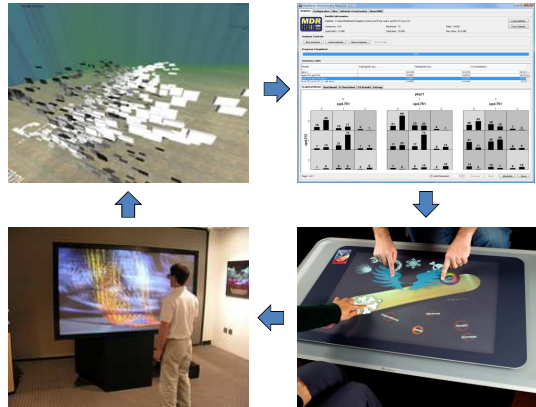
Scientific Visualization



Information Visualization



Visual Analytics



EDUCATION

Introductory Science and Mathematics Education for 21st-Century Biologists

William Bialek^{1,3} and David Botstein^{2,3*}

Galileo wrote that "the book of nature is written in the language of mathematics"; his quantitative approach to understanding the natural world arguably marks the beginning of modern science. Nearly 400 years later, the fragmented teaching of science in our universities still leaves biology outside the quantitative and mathematical culture that has come to define the physical sciences and engineering. This strikes us as particularly inopportune at a time when opportunities for quantitative thinking about biological systems are exploding. We propose that a way out of this dilemma is a unified introductory science curriculum that fully incorporates mathematics and quantitative thinking.

These traditions have resulted in a deep bifurcation in culture and quantitative competence among the scientific disciplines. On one branch are mathematics, the physical sciences, and engineering. Scientists educated along this branch achieve a high level of quantitative expertise: They generally have some mastery over and comfort with not only multivariate calculus and differential equa-



- chgr.mc.vanderbilt.edu/bushlab
- www.gettinggeneticsdone.com
- william.s.bush@vanderbilt.edu