# Bioinformatics

### Jason H. Moore, Ph.D.

Frank Lane Research Scholar in Computational Genetics
Associate Professor of Genetics
Adjunct Associate Professor of Biological Sciences
Adjunct Associate Professor of Community and Family Medicine
Dartmouth College

Affiliate Associate Professor of Computer Science
University of New Hampshire

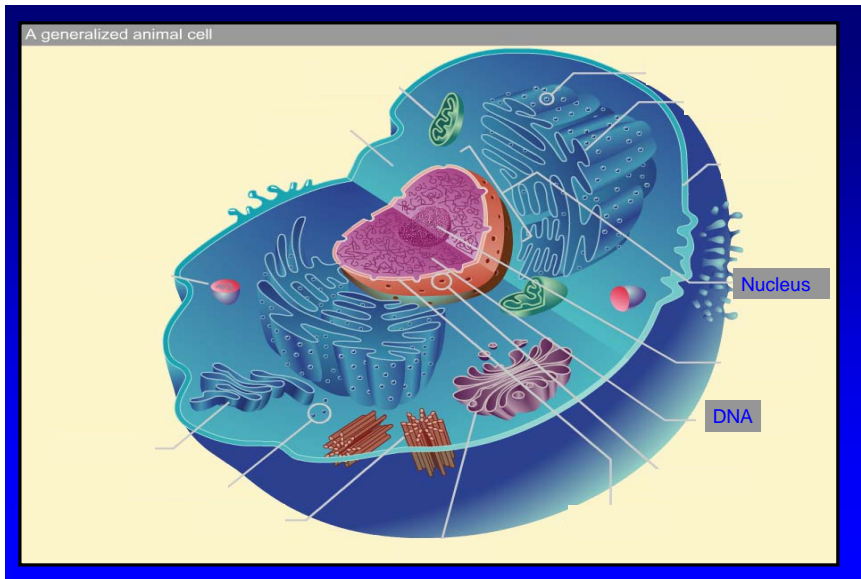Adjunct Associate Professor of Computer Science
University of Vermont

www.epistasis.org
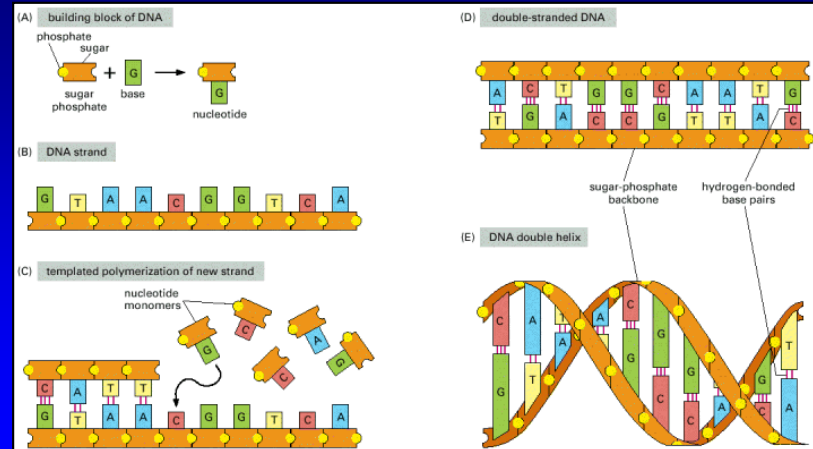jason.h.moore@dartmouth.edu
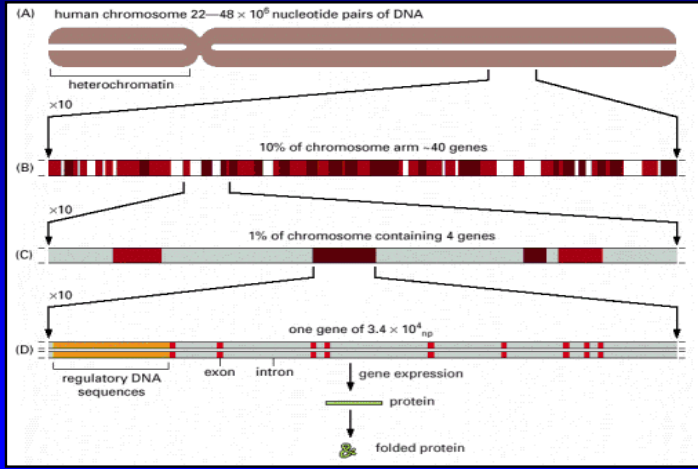
1

© Jason H. Moore

---

# Genotype -> Phenotype

2

---



A generalized animal cell

Nucleus

DNA

---

# DNA



(A) building block of DNA
phosphate
sugar
sugar   base
phosphate
nucleotide

(B) DNA strand

(C) templated polymerization of new strand
nucleotide monomers

(D) double-stranded DNA

sugar-phosphate backbone
hydrogen-bonded base pairs

(E) DNA double helix

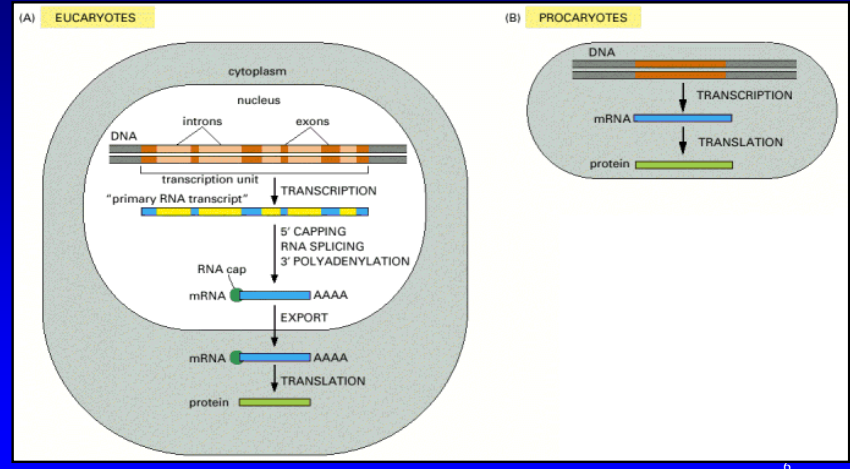Alberts et al. 2002

4

---

# Chromosome Structure
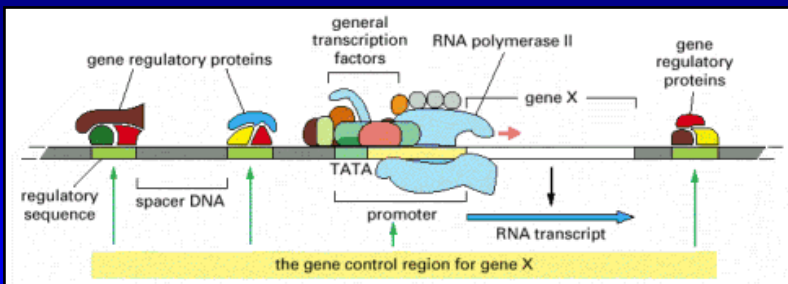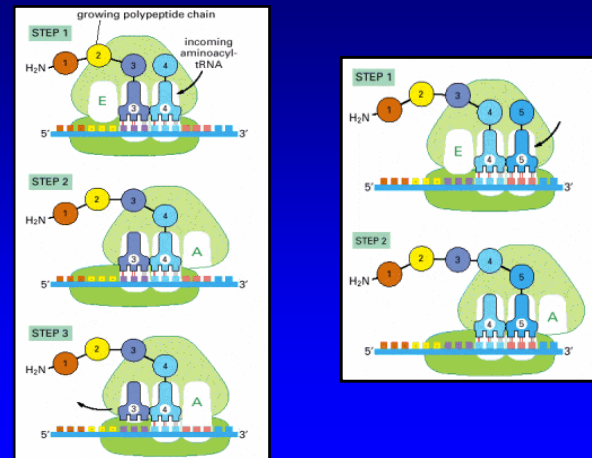


Alberts et al. 2002

# Transcription



Alberts et al. 2002

# Regulation of Transcription



Alberts et al. 2002

# Translation
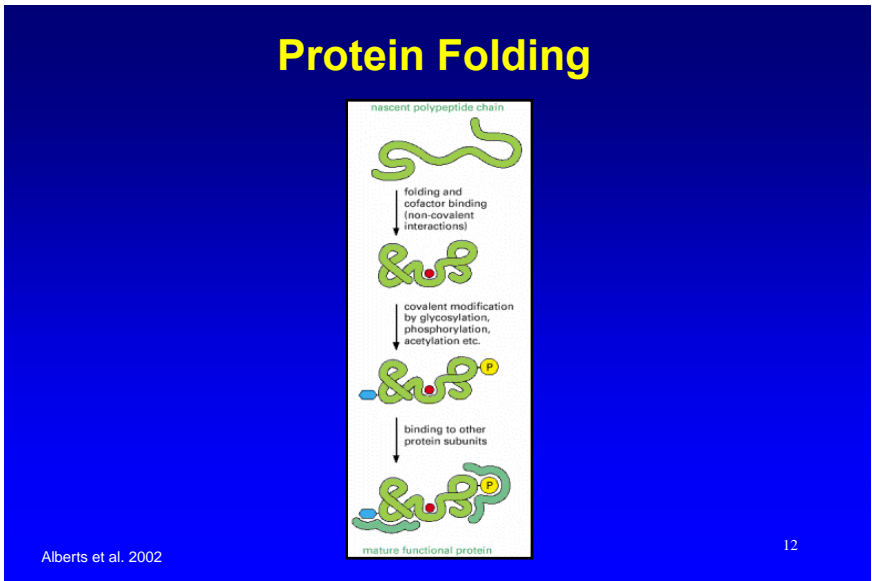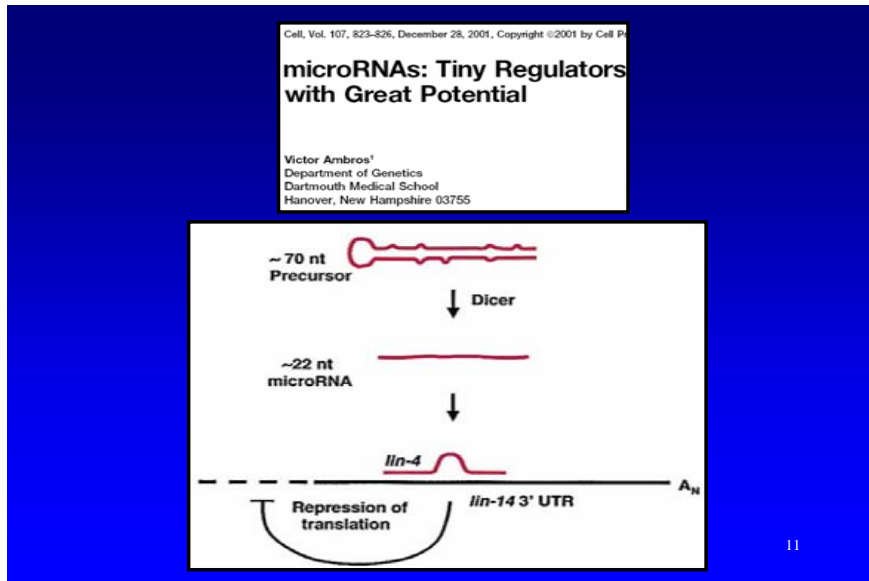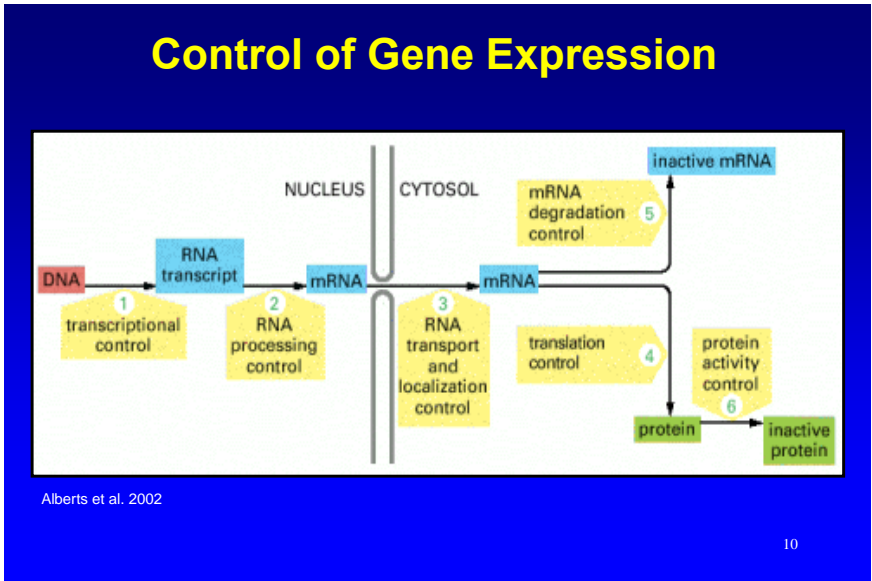


Alberts et al. 2002

Alberts et al. 2002

# Control of Gene Expression



Alberts et al. 2002

10



11

# Protein Folding



Alberts et al. 2002

12

## Amino Acids

| 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 | | | |

Alberts et al. 2002

13

## Amino Acids → Peptides → Proteins



Brown 2002

14



Copyright ©1993 The Scripps Research Institute

Water  Amino acids  *Alanine*  *Tryptophan*  DNA helix  Protein helix  Enzyme  *Cu, Zn Superoxide Dismutase*



Copyright ©1993 The Scripps Research Institute

Antibody  *Dob model*  Virus  *Reovirus core*

3438

Copyright ©1993 The Scripps Research Institute

# Gene Networks



18

# Protein Networks

Barabasi, *Scientific American* (2003)



© Jason H. Moore

19

# Genetic Architecture



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

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

20

3439

## The Tree of Life



http://tolweb.org/tree/

© Jason H. Moore

21

## Molecular Phylogenetics



Brown 2002

22

## Measuring DNA

23



July 3, 2000

24

3440

NATURE REVIEWS | **GENETICS**
VOLUME 6 | APRIL 2005 | **333**

OPINION

The Human Genome Diversity
Project: past, present and future

*L. Luca Cavalli-Sforza*

**Africans**
1 Bantu
2 Mandenka
3 Yoruba
4 San
5 Mbuti pygmy
6 Biaka
7 Mozabite

**Europeans**
8 Orcadian
9 Adygei
10 Russian
11 Basque
12 French
13 North Italian
14 Sardinian
15 Tuscan

**Central and Southern Asians**
19 Balochi
20 Brahui
21 Makrani
22 Sindhi
23 Pathan
24 Burusho
25 Hazara
26 Uygur
27 Kalash

**Western Asians**
16 Bedouin
17 Druze
18 Palestinian

**Eastern Asians**
28 Han (S. China)
29 Han (N. China)
30 Dai
31 Daur
32 Hezhen
33 Lahu
34 Miao
35 Oroqen
36 She
37 Tujia
38 Tu
39 Xibo
40 Yi
41 Mongola
42 Naxi
43 Cambodian
44 Japanese
45 Yakut

**Oceanians**
46 Melanesian
47 Papuan

**Native Americans**
48 Karitiana
49 Surui
50 Colombian
51 Maya
52 Pima

25

## Transposon Repeat Polymorphism

**Sequence variation in the human angiotensin converting enzyme**

nature genetics • volume 22 • may 1999

Mark J. Rieder[1], Scott L. Taylor[1], Andrew G. Clark[2] & Deborah A. Nickerson[1]

```
1   ggctgggcgt ggtggctcaa gcctgtaatc ccagcacttt gggaggctga ggtgggcgca
61  tcgcttgagc ccaggagttc aagaccagcc tggccaacat cgcaaaacct cgtctctaca
121 aaaaaaaaat agctgggctt ggtggtgcgt gcacctacag tcccagctac tcttgaaact
181 gaggggggaag gatcacctga gcccaggagg tcaaggctac agtgagctgt gattgcacta
241 ctgcaccccca gcctgcgtga cagagtgaga cctccccccca aaaaaaagag agagagaaaa
```

26

## SNPs
### Single Nucleotide Polymorphisms

| Subject #1 | Subject #2 | Subject #3 |
|---|---|---|
| -- AG**G**TCA -- | -- AG**G**TCA -- | -- AG**C**TCA -- |
| -- AG**G**TCA -- | -- AG**C**TCA -- | -- AG**C**TCA -- |

Two *alleles* (**G** and **C**)

Three *genotypes* (**GG**, **GC**, **CC**)

27

## SNPs
### Single Nucleotide Polymorphisms

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

28

**AFFYMETRIX**

## Data Sheet

GeneChip® Human Mapping 500K Array Set

**Figure 1:** GeneChip® Mapping Assay Overview.

Genomic DNA (250 ng) — RE Digestion

Xba    Xba    Xba

Adapter Ligation

PCR: One Primer Amplification

Complexity Reduction

Fragmentation and End-Labeling

Hybridization & Wash

AA BB AB

# Measuring RNA

30

# cDNA Microarrays



Reference RNA (pooled from cell lines)

Tumor RNA

RNA reverse transcribed with Cy3 or Cy5 labeled nucleotides

Reference sample cDNA

Tumor cDNA

cDNAs mixed and hybridized to microarray

ERBB2

Mol Interv. 2002    image of scanned microarray

31

## APPLICATIONS OF DNA MICROARRAYS IN BIOLOGY

Roland B. Stoughton
*GHC Technologies, Incorporated, La Jolla, California 92037;*
*email: roland_stoughton@ghctechnologies.com*

32

3442

## Process Flow for Microarrays



## 2D Gels and Mass Spectrometry



Metabolic Engineering **4**, 98–106 (2002)

## Measuring Proteins

# Protein Profiling in Tissues



Am J Pathol. 2004

37

# Protein Profiling in Tissues



Am J Pathol. 2004

# Tissue Microarrays



Nat Rev Drug Discov. 2003   39

# Emerging Technologies

40

## Nanotechnology



Current Opinion in Biotechnology

41

## Nanotechnology



Nature Reviews Cancer 2005

42

## Lab-on-a-Chip



Current Opinion in Structural Biology 2003

43

## Databases

44

# Bioinformatics: Databases



101100 010010

$\Sigma$

45

# http://www.ncbi.nlm.nih.gov



46

# http://www.ncbi.nlm.nih.gov/gquery/gquery.fcgi



47

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



48

3446

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

www.pharmgkb.org

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

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

49

50

51

52

# Analysis

53

# 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?

54

# Objectives

Data $\longrightarrow$ Variable Selection $\longrightarrow$ Statistical Modeling $\longrightarrow$ Prediction Classification

$I = a_1 x_1 + a_2 x_2$ $\longrightarrow$ Tumor A

Tumor B

Iterate

55

# Objectives

Data $\longrightarrow$ Variable Selection $\longrightarrow$ Statistical Modeling $\longrightarrow$ Prediction Classification

$I = a_1 x_1 + a_2 x_2$ $\longrightarrow$ Tumor A

Tumor B

Iterate

56

## Hypothesis Testing
*"The truth is out there"*

**Truth**

| | | $H_o$ False | $H_o$ True |
|---|---|---|---|
| Decision | Reject $H_o$ | Yes! | Type I Error |
| | Accept $H_o$ | Type II Error | Yes! |

**Type I error**
**Type II error**
**Type III error**

57

© Jason H. Moore

## Cross-Validation (CV)

- Data-driven methods susceptible to overfitting
- Biological datasets often have more variables than observations (i.e. wide data)
- The value of any statistical model is its ability to make predictions in new data
- Cross-validation (CV) allows generalizability to be estimated in a single dataset

58

© Jason H. Moore

## Cross-Validation (CV)

- CV uses independent portions of the data to estimate the testing accuracy



Train 9/10  Train 9/10  Train 9/10
Test 1/10   Test 1/10   Test 1/10

59

© Jason H. Moore

## Cross-Validation (CV)
Hastie et al. *The Elements of Statistical Learning* (2001)
Ripley BD. *Pattern Recognition and Neural Networks* (1996)

- Leave One Out Cross Validation (LOOCV)
  - Better for small datasets
  - Unbiased estimate of prediction error
  - High variance due to similarity of training sets

- *n*-fold CV (e.g. 5-fold or 10-fold CV)
  - Better for larger datasets
  - Estimate of prediction error may be biased
  - Lower variance
  - May need to repeat several times and average results

60

© Jason H. Moore

3449

## Cross-Validation Consistency (CVC)

Ritchie et al., *American Journal of Human Genetics* 69:138-147 (2001)
Moore et al., *Genetic Epidemiology* 23:57-69 (2002)
Moore, *Lecture Notes in Computer Science* 2611, Springer-Verlag, Berlin (2003).

- CV can be difficult with data-driven methods
- Can find different models with each CV dataset
- CVC is a measure of how consistently particular variables or combination of variables are identified in each CV interval.
- CVC can be used as a measure of association
- Once important variables are found, a model can be fit to the entire dataset using just those variables.

61

© Jason H. Moore

## Cross-Validation Consistency (CVC)

Ritchie et al., *American Journal of Human Genetics* 69:138-147 (2001)
Moore et al., *Genetic Epidemiology* 23:57-69 (2002)
Moore, *Lecture Notes in Computer Science* 2611, Springer-Verlag, Berlin (2003).

### CVC Example with 5-fold CV

Data Interval        Genes Identified in Best Model

*A,C*,E

*A,C*,H                    Significant Genes

*A*,B,*C*,F      ➜      *A and C*

*A,C*,D

*A*,G,H

62

© Jason H. Moore

## Generalizability: The Three-Way Data Split

Rowland, *Lecture Notes in Computer Science* 2611, Springer-Verlag, Berlin (2003).

Train     Test | Generalize  Train | Test  Generalize

Generalize | Test | Train

1. Choose model that min[ abs( $E_{train} - E_{test}$ ) ]
2. Evaluate generalizability of the final model
3. The model with the best prediction error may not generalize the best.

63

© Jason H. Moore

## 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.

64

© Jason H. Moore

3450

## Permutation Testing

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

### Example

| Y | X1 X2 |
|---|-------|
| 1 | AA Bb |
| 1 | Aa BB |
| 1 | aa Bb |
| 0 | AA Bb |
| 0 | Aa bb |
| 0 | aa bb |

Permute Y →

| Y | X1 X2 |
|---|-------|
| 0 | AA Bb |
| 1 | Aa BB |
| 0 | aa Bb |
| 0 | AA Bb |
| 1 | Aa bb |
| 1 | aa bb |

Calculate Statistic → $Z$

Repeat ←

65

© Jason H. Moore

## 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

Critical region, $\alpha$

$Z$

66

© Jason H. Moore

## 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

Population

Random Sample

Inference

### Distribution Unknown

Population

Random Sample

Random Samples

Inference

67

© Jason H. Moore

## Genetic Programming

68

Is GP an appropriate computational tool for solving complex biological problems??

69

Vanilla GP - NO!

GP

70

Modern GP - YES!

Post-Processing
Parameter Sweep
Pre-Processing
Expert Knowledge
GP

71

Case Study:
Symbolic Modeling

72

3452

## Genetic Architecture



POTENTIAL REACTION SURFACE

ENVIRONMENTAL CHANGE

STENOSIS OF CORONARY ARTERIES

POSSIBLE ENVIRONMENTS

?

AGE

EXERCISE

DIET

LIPID METABOLISM

HAEMOSTASIS

CARBOHYDRATE METABOLISM

BLOOD PRESSURE REGULATION

STRESS

SMOKING

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

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

73

---

## Genetic Analysis of Atrial Fibrillation

Tsai et al., *Circulation* (2004)
Moore et al., *Journal of Theoretical Biology* (2006)

- 250 consecutive patients admitted with AF from Taipei, Taiwan
- 250 age and gender matched controls
- 3 candidate genes
  - *Angiotensin converting enzyme (ACE)* gene
    - I/D polymorphism
  - *Angotensinogen* (*AGT*)
    - *T174M, M235T, G-6A, A-20C, G-152A*, and *G-217A* polymorphisms
  - *Angiotensin II type I receptor* (*AT-1*)
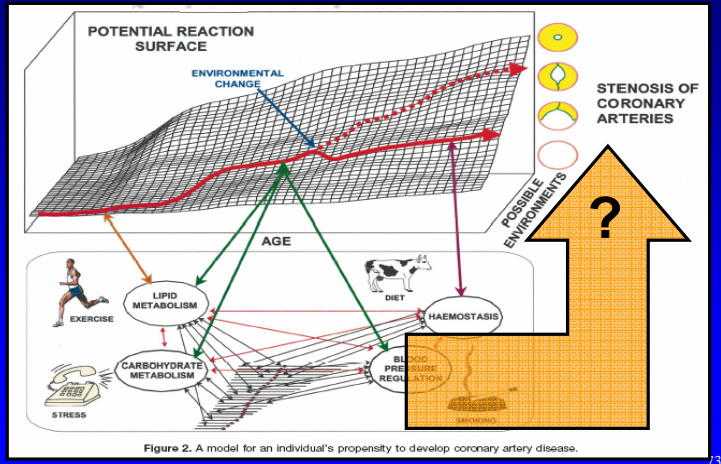    - *A1166C* polymorphism

74

---

## Renin-Angiotensin System

Tsai et al., *Circulation* (2004)
Moore et al., *Journal of Theoretical Biology* (2006)



Renin signal

Renin genes

ACE genes

V

V  V  ACE inhibition

V  AGT genes

Renin

ACE

Bradykinin

AGT

ANG 1

ANG II

BP

Flow factor
Variable
Flow Control
Gene product / intermediate
Flow

75

---

## Symbolic Discriminant Analysis

Moore et al., In: De Raedt, L., Flach, P. (eds) *Lecture Notes in Artificial Intelligence* 2167 (2001)
Moore et al., *Genetic Epidemiology* 23, 57-69 (2002)
Moore, *Lecture Notes in Computer Science* 2611, Springer-Verlag, Berlin (2003)

- Supervised classification approach
- Can use GP to build discriminant functions
- Accuracy is fitness function



$Y = X_1 + X_2 * X_3$

Frequency

A    B

Symbolic Discriminant Scores

76

3453

## SDA Modeling using Modern GP
Moore et al., *Human Heredity* (2007)

- Parameter Sweeps
  - STEP 1: Full factorial experimental design
  - STEP 2: ANOVA
- Coarse-Grained Search
  - STEP 3: $10^5$ SDA runs
- Expert Knowledge
  - STEP 4: Statistical model of 100 best trees
- Fine-Grained Knowledge-Based Search
  - STEP 5: $10^9$ SDA runs using EDA
- Model Interpretation
  - STEP 6: Function mapping
  - STEP 7: Interaction dendrogram

77

© Jason H. Moore

---

## Cross-Validation Strategy
Moore et al., *Human Heredity,* (2007)



78

© Jason H. Moore

---

## Parameter Sweep
Moore et al., *Human Heredity,* (2007)

- STEP 1: Full factorial experimental design
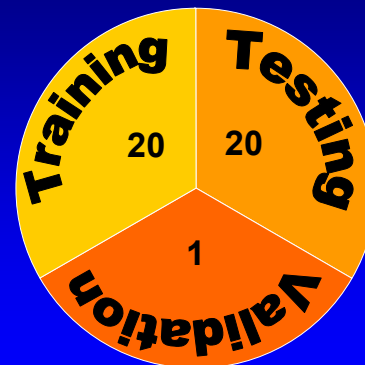  - Population Size (P): {100, 500, 1000}
  - Generations (G): {100, 500, 1000}
  - Tree Depth (T): {1, 2, 3}
  - Function Set (F):
    - 1: {+, -, *, /}
    - 2: {<, >, <=, >=, =, !=, max, min}
    - 3: 12
    - 4: {AND, OR, NOT, NOR, XOR}
    - 5: 14
    - 6: 24
    - 7: 124
  - 3*3*3*7 = 189 level combinations
  - 189 levels * 10 random seeds = 1890 SDA runs

79

© Jason H. Moore

---

## Parameter Sweep
Moore et al., *Human Heredity,* (2007)

- STEP 2: 4-Way Analysis of Variance (ANOVA)

| Treatment | DF | F | P-Value |
|-----------|-----|--------|---------|
| Depth | 2 | 715.6 | <0.001 |
| Function | 6 | 1028.2 | <0.001 |
| Generations | 2 | 1.8 | 0.161 |
| Population | 2 | 10.7 | <0.001 |
| D*F*P | 24 | 2.4 | 0.001 |

80

© Jason H. Moore

3454

## Parameter Sweep
Moore et al., *Human Heredity,* (2007)



P = 0.001

© Jason H. Moore

## Course-Grained Search
Moore et al., *Human Heredity,* (2007)

- STEP 3: $10^5$ SDA runs

| Population Size | 500 |
|---|---|
| Generations | 100 |
| Tree Depth | 3 |
| Crossover | Single-point subtree |
| Crossover Frequency | 0.9 |
| Mutation Frequency | 0.01 |
| Fitness Function | Accuracy |
| Selection | Three-way Tournament |
| Function Set | +, -, *, /, =, !=, <, >, ≤, ≥, min, max |
| Terminal Set | -2, -1, 0, 1, 2, M235T, T174M, G6A, A20C, G152A, G217A, ACE I/D, AT1R |

© Jason H. Moore

## Generating Expert Knowledge
Moore et al., *Human Heredity,* (2007)

- STEP 4: Statistical modeling of 100 best trees



© Jason H. Moore

## Fine-Grained Search
Moore et al., *Human Heredity,* (2007)

- STEP 5: $10^9$ SDA runs using an EDA



**Accuracy = 0.644**

© Jason H. Moore

3455

## Fine-Grained Search
Moore et al., *Human Heredity,* (2007)

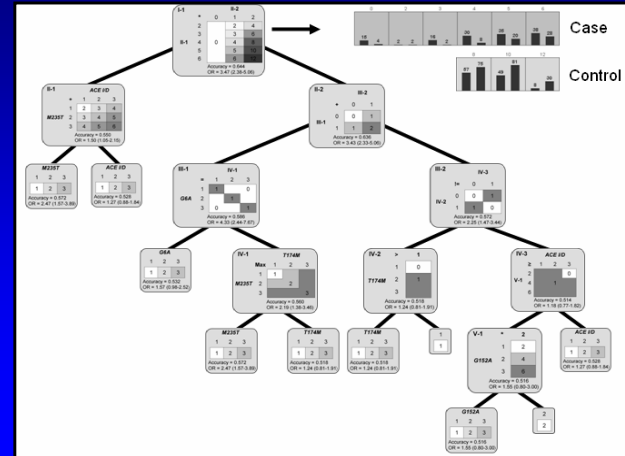- STEP 5: $10^9$ SDA runs using an EDA



**Accuracy = 0.644**

© Jason H. Moore

85

## Interpretation
Moore et al., *Human Heredity,* (2007)
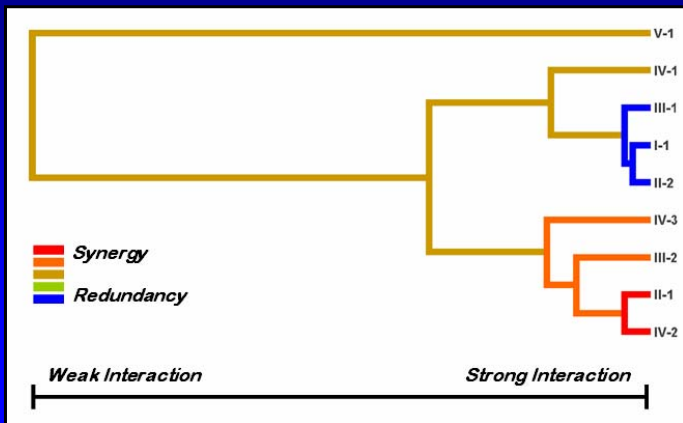
- STEP 6: Function Mapping



© Jason H. Moore

86

## Interpretation
Moore et al., *Human Heredity,* (2007)
Moore et al., *Journal of Theoretical Biology* (2006)
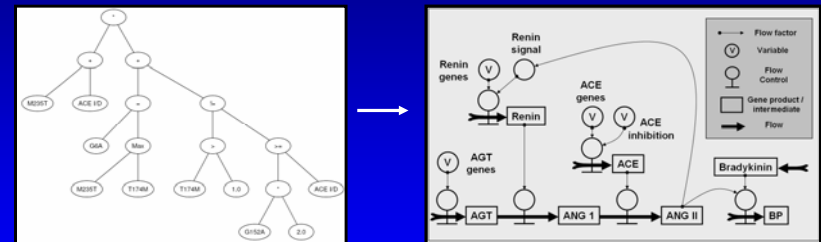
- STEP 6: Interaction Dendrogram



*Synergy*

*Redundancy*

*Weak Interaction*          *Strong Interaction*

© Jason H. Moore

87

## Interpretation
Moore et al., *Human Heredity,* (2007)



© Jason H. Moore

88

3456

## Interpretation

Moore, *Nature Genetics* (2005)
Moore and Williams, *BioEssays* (2005)



© Jason H. Moore

89

## Modern GP - YES!



90

## Additional Examples

- Moore, J.H., White, B.C. Genome-wide genetic analysis using genetic programming: The critical need for expert knowledge. In: Genetic Programming Theory and Practice IV, in press, Springer (2006).

- Moore, J.H., White, B.C. Exploiting expert knowledge in genetic programming for genome-wide genetic analysis. Lecture Notes in Computer Science, in press (2006).

- Moore, J.H. Genome-wide analysis of epistasis using multifactor dimensionality reduction: feature selection and construction in the domain of human genetics. In: Knowledge Discovery and Data Mining: Challenges and Realities with Real World Data, IGI, in press (2006).

© Jason H. Moore

91

## Acknowledgments

**Nate Barney**
**Bill White**

**NIH R01 AI59694, LM009012**

**www.epistasis.org**
**compgen.blogspot.com**
**jason.h.moore@dartmouth.edu**

© Jason H. Moore

92