# http://tinyurl.com/666-week0

Welcome to Biostatistics 666!

Please fill in the survey while we wait to start.

# Course Overview and Welcome!

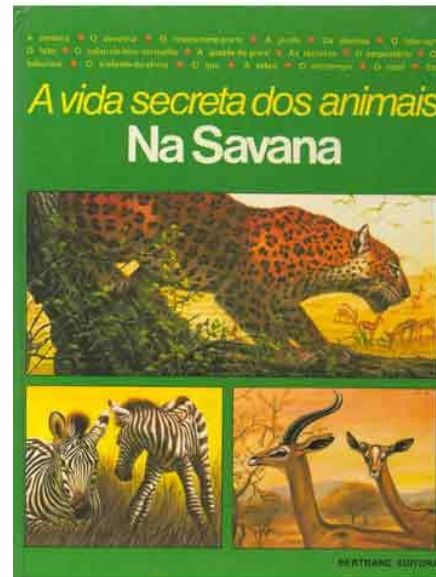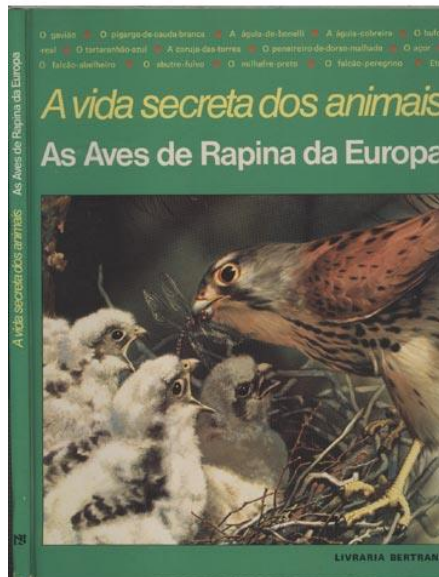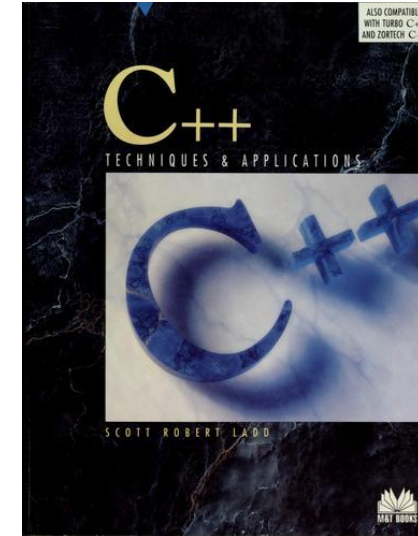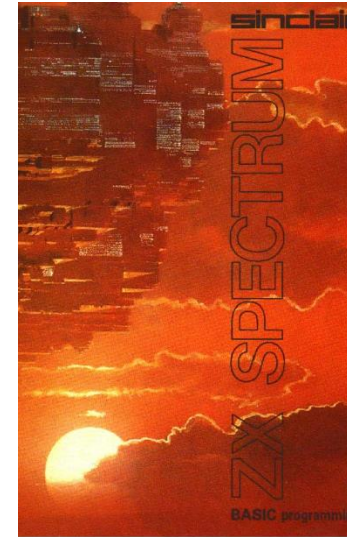Biostatistics 666

Goncalo Abecasis

# My Day Job…

# Use Genetic Variation, Math and Computation to Understand Human Disease

# Why?



```
5 REM pangolins
10 LET nq=100: REM number of questions and animals
15 DIM q$(nq,50): DIM a(nq,2): DIM r$(1)
20 LET qf=8
30 FOR n=1 TO qf/2−1
40 READ q$(n): READ a(n,1): READ a(n,2)
50 NEXT n
60 FOR n=n TO qf−1
70 READ q$(n): NEXT n

100 REM start playing
110 PRINT "Think of an animal.","Press any key to continue."
120 PAUSE 0
130 LET c=1: REM start with 1st question
140 IF a(c,1)=0 THEN GO TO 300
150 LET p$=q$(c): GO SUB 910
160 PRINT "?": GO SUB 1000
170 LET in=1: IF r$="y" THEN GO TO 210
180 IF r$="Y" THEN GO TO 210
190 LET in=2: IF r$="n" THEN GO TO 210
200 IF r$<>"N" THEN GO TO 150
210 LET c=a(c,in): GO TO 140

300 REM animal
310 PRINT "Are you thinking of"
320 LET p$=q$(c): GO SUB 900: PRINT "?"
330 GO SUB 1000
340 IF r$="y" THEN GO TO 400
350 IF r$="Y" THEN GO TO 400
360 IF r$="n" THEN GO TO 500
370 IF r$="N" THEN GO TO 500
```
211

# Human Genetics, Sample Sizes over My Time

| Year | No. of Samples | No. of Markers | Publication |
|------|----------------|----------------|-------------|
| Ongoing | 120,000 | 600 million | NHLBI Precision Medicine Cohorts / TopMed |
| 2016 | 32,488 | 40 million | Haplotype Reference Consortium (Nature Genetics) |
| 2015 | 2,500 | 80 million | The 1000 Genomes Project (Nature) |
| 2012 | 1,092 | 40 million | The 1000 Genomes Project (Nature) |
| 2010 | 179 | 16 million | The 1000 Genomes Project (Nature) |
| 2010 | 100,184 | 2.5 million | Lipid GWAS (Nature) |
| 2008 | 8,816 | 2.5 million | Lipid GWAS (Nature Genetics) |
| 2007 | 270 | 3.1 million | HapMap (Nature) |
| 2005 | 270 | 1 million | HapMap (Nature) |
| 2003 | 80 | 10,000 | Chr. 19 Variation Map (Nature Genetics) |
| 2002 | 218 | 1,500 | Chr. 22 Variation Map (Nature) |
| 2001 | 800 | 127 | Three Region Variation Map (Am J Hum Genet) |
| 2000 | 820 | 26 | T-cell receptor variation (Hum Mol Genet) |

# Course Logistics

Scheduling

Office Hours

Class Notes

Grading

# Course Objective

- Introduce statistical models used in gene mapping studies

- Survey common algorithms used for handling genetic data

- Provide foundation for using gene mapping methods

- Provide foundation for refining and developing gene mapping methods

# Course Notes

- We will not be using a textbook
  - Extremely important to attend class, and ask questions as needed!

- Copies of slides and additional content available online at
  - http://genome.sph.umich.edu/wiki/666

# Assessment

- Grading will be a combination of:
  - Home work assignments (40%, approximately weekly)
  - In-class written assessments (60%, two of these)

- Greg Zajac (Biostatistics PhD student) will be helping me with grading.

# Academic Integrity

- All assignments you submit for evaluation must represent your own work.

- If you copy or paraphrase others, you must clearly mark these sections and indicate sources.

- Any involvement in cheating constitutes academic misconduct and is a serious offense.

- See also the School policy on academic conduct.

# Academic Integrity

- All assignments you submit for evaluation must represent your own work.

- If you copy or paraphrase others, you must clearly mark these sections and indicate sources.

  If a set of assignments or exams is broadly identical, each will be scored as zero and referred to Department or School.

- Any involvement in cheating constitutes academic misconduct and is a serious offense.

- See also the School policy on academic conduct.

# Scheduling

- We will try to start classes at 8:30 sharp.

  - Due to prior commitments, I may have to miss several lectures and starting sharply on time should allow us to make up for any lost time.

# Office Hours

- We will try to find a time …

- To provide input, please fill in times when **you are available** at:

    **http://tinyurl.com/666-office-hours**

- **Thanks!**

# Goals for Today …

- Overview of the evolution of complex disease studies
  - Current state of the art, challenges, opportunities

- Heritability
  - Estimating the total (additive) contribution of genetic variation to a trait

- Genomewide Association Studies
  - Linkage Disequilibrium, Genotype imputation

- Sequence-based Association Studies

# If you are new to genetics …

- 23andMe provides a nice intro to the basic principles of genetics

- http://www.23andme.com/gen101/

- The material is all quite good and the videos are easy to watch.

- If genetics are new to you, I recommend you browse the first 4 sections.

# Modern Gene Mapping Studies

A quick overview!

# How human genetic studies work …

- DNA is our instruction manual

- We are all built mostly to the same plan…
  - Any two human DNA molecules are ~99.9% the same

- We each have our manual, with small variations from the typical plan
  - Most of these variations, as far as know, have no health consequences
  - Some modify key cell processes, a lot (these are very rare) or a little (more commonly)

- Search for variants that modify interesting health outcomes…
- … identify the cellular processes they modify …
- … improve understanding of biology, identify therapeutic targets …

# Human Genetics, Sample Sizes over My Time

| Year | No. of Samples | No. of Markers | Publication |
|---|---|---|---|
| Ongoing | 120,000 | 600 million | NHLBI Precision Medicine Cohorts / TopMed |
| 2016 | 32,488 | 40 million | Haplotype Reference Consortium (Nature Genetics) |
| 2015 | 2,500 | 80 million | The 1000 Genomes Project (Nature) |
| 2012 | 1,092 | 40 million | The 1000 Genomes Project (Nature) |
| 2010 | 179 | 16 million | The 1000 Genomes Project (Nature) |
| 2010 | 100,184 | 2.5 million | Lipid GWAS (Nature) |
| 2008 | 8,816 | 2.5 million | Lipid GWAS (Nature Genetics) |
| 2007 | 270 | 3.1 million | HapMap (Nature) |
| 2005 | 270 | 1 million | HapMap (Nature) |
| 2003 | 80 | 10,000 | Chr. 19 Variation Map (Nature Genetics) |
| 2002 | 218 | 1,500 | Chr. 22 Variation Map (Nature) |
| 2001 | 800 | 127 | Three Region Variation Map (Am J Hum Genet) |
| 2000 | 820 | 26 | T-cell receptor variation (Hum Mol Genet) |

# A comprehensive review of genetic association studies

Joel N. Hirschhorn, MD, PhD[1-3] , Kirk Lohmueller[1], Edward Byrne[1], and Kurt Hirschhorn, MD[4]

"... of the 166 associations which have been studied 3 or more times, only six have been consistently replicated."

Hirschhorn et al (2002)

# A Genomewide Study of Obesity



Seven of eight confirmed BMI loci show strongest expression in the brain…

Willer et al, *Nature Genetics,* 2009

# Current State of Genetic Association Studies

- Surveying common variation across 10,000s - 100,000s of individuals is now routine, using genotyping arrays

- Many common alleles have been associated with a variety of human complex traits

- The functional consequences of these alleles are often subtle, and translating the results into mechanistic insights remains challenging

- Sequencing studies are starting to allow studies to extend to rare variants, which can lead to easier to understand biology

# Heritability

How Much of Phenotypic Variation Can Genetic Variation Explain?

Does Genetic Similarity Predict Phenotypic Similarity?

# Variability in Height



We might often summarize this distribution with a mean and variance.

# Variability in Height, Pairs of Observations

If sampling pairs of individuals, we might also record covariance between pairs of observations …



(Data from David Duffy)

# Height in DZ and MZ twins



(How would you interpret these data from David Duffy?)

# Variance-Covariance Matrix

$$\Omega = \begin{bmatrix} V(y_1) & Cov(y_1, y_2) \\ Cov(y_1, y_2) & V(y_2) \end{bmatrix}$$

Model describes not only variance of each
observation but also covariance for pairs of observations

# A Simple Model for the Variance-Covariance Matrix

$$\Omega = \begin{bmatrix} \sigma_g^2 + \sigma_e^2 & 2\varphi\sigma_g^2 \\ 2\varphi\sigma_g^2 & \sigma_g^2 + \sigma_e^2 \end{bmatrix}$$

Where,

$\varphi$ is the kinship coefficient for the two individuals

# Linkage Disequilibrium and Genetic Association Studies

Genetic Association Signals at Nearby Variants …

# Linkage Disequilibrium

- Chromosomes are mosaics

- Extent and conservation of mosaic pieces depends on
  - Recombination rate
  - Mutation rate
  - Population size
  - Natural selection

- Combinations of alleles at very close markers reflect ancestral haplotypes

**Ancestor**

**Present-day**

# Obesity and the *NEGR1* locus



Multiple nearby SNPs show evidence for association with obesity.
The associated alleles usually appear together, in a haplotype.

Willer et al, *Nature Genetics, 2009*

# Observed Genotypes

**Observed Genotypes**

**Study Sample**
Inexpensive measurements
at 100,000s of markers

. . . . . A . . . . . . . . A . . . . . A . . .
. . . . G . . . . . . . . C . . . . A . . .

**Reference Haplotypes**

**Reference Sample**
Detailed measurements
of 1,000,000s of markers

C G A G A T C T C C T T C T T C T G T G C
C G A G A T C T C C C G A C C T C A T G G
C C A A G C T C T T T T C T T C T G T G C
C G A A G C T C T T T T C T T C T G T G C
C G A G A C T C T C C G A C C T T A T G C
T G G G A T C T C C C G A C C T C A T G G
C G A G A T C T C C C G A C C T T G T G C
C G A G A C T C T T T T C T T T T G T A C
C G A G A C T C T C C G A C C T C G T G C
C G A A G C T C T T T T C T T C T G T G C

# Identify Match Among Reference

# Fill-in Missing Genotypes

**Observed Genotypes**

c g a g A t c t c c c g A c c t c A t g g
c g a a G c t c t t t t C t t t c A t g g

**Reference Haplotypes**

C G A G A T C T C C T T C T T C T G T G C
C G A G A T C T C C C G A C C T C A T G G
C C A A G C T C T T T T C T T C T G T G C
C G A A G C T C T T T T C T T C T G T G C
C G A G A C T C T C C G A C C T T A T G C
T G G G A T C T C C C G A C C T C A T G G
C G A G A T C T C C C G A C C T T G T G C
C G A G A C T C T T T T C T T T T G T A C
C G A G A C T C T C C G A C C T C G T G C
C G A A G C T C T T T T C T T C T G T G C

# The Role of Sequencing
# in Genetic Association Studies

# Shotgun Sequence Data



TAGCTGATAGCTAGATAGCTGATGAGCCCGAT

ATAGCTAGATAGCTGATGAGCCCGATCGCTGCTAGCTC

ATGCTAGCTGATAGCTAGCTAGCTGATGAGCC

AGCTGATAGCTAGCTAGCTGATGAGCCCGATCGCTG

GCTAGCTGATAGCTAGCTAGCTGATGAGCCCGA

Sequence Reads

5'-ACTGGTCGATGCTAGCTGATAGCTAGCTAGCTGATGAGCCCGATCGCTGCTAGCTCGACG-3'

Reference Genome

**A/C**

Predicted Genotype

The 1000 Genomes Project (2008 – 2015)

# Optimal Model for Analyzing 1000 Genomes?

| 1000 Genomes Call Set (CEU) | Homozygous Reference Error | Heterozygote Error | Homozygous Non-Reference Error |
|---|---|---|---|
| Broad | 0.66 | 4.29 | 3.80 |
| Michigan | 0.68 | 3.26 | 3.06 |
| Sanger | 1.27 | 3.43 | 2.60 |

- Michigan caller combines …
  - Markov models to identify shared haplotypes,
  - Classifiers to distinguish true variants from error,
  - Strategies to distribute computation across cluster

# Optimal Model for Analyzing 1000 Genomes?

| 1000 Genomes Call Set (CEU) | Homozygous Reference Error | Heterozygote Error | Homozygous Non-Reference Error |
|---|---|---|---|
| Broad | 0.66 | 4.29 | 3.80 |
| Michigan | 0.68 | 3.26 | 3.06 |
| Sanger | 1.27 | 3.43 | 2.60 |
| Majority Consensus | 0.45 | 2.05 | 2.21 |

- Common to see **"ensemble" methods outperform the best single method**

# A Key Goal of
# Sequence Based Association Studies

**UNDERSTAND FUNCTION**

**LINKING EACH LOCUS TO DISEASE**

**What happens in gene knockouts?**

- Use sequencing to find rare human "knockout" alleles
- Why? Results of animal studies an *in vitro* studies often murky
- The challenge? Natural knockouts are extremely rare

# TOPMed Sequencing as of May 25, 2017

- 76,436 genomes
    - 74,890 pass quality checks          (98.0%)
    - 946 flagged for low coverage    (  1.2%)
    - 606 fail quality checks          (  0.8%)


- Mean depth:        38.3x

- Genome covered: 98.6%

- Contamination:     0.28%


- $10^{16}$ sequenced bases

**Overall Genome Counts**

● Pass    ● Flag    ● Fail

67,317

# $10^{16}$ sequenced bases



Number of snowflakes covering ~10 square miles in a 10-inch deep snowstorm.
100x more data than the 1000 Genomes Project.

# $10^{16}$ sequenced bases



US corn production in 2014: $1.3 \times 10^{15}$ kernels

# Browse All Variations Online
## http://bravo.sph.umich.edu

Peter VandeHaar

## KMT2D



## PCSK9



496 missense, 26 inframe indels, 0 stop or frameshifts

91 missense, 4 inframe indels, 7 stop or frameshifts

# How to help TOPMed advance discoveries?

- Genomewide analyses at scale are challenging

- Even simple analysis can require 1,000s of CPU days to complete

- Need to engage diverse teams in analysis and interpretation



```
snp,pvalue
rs1234,0.05
rs4343,0.0002
rs51101,0.61
rs981,0.000018
rs2223,0.72
```

- Exploring new ways to engage populations in research

- Continuous Engagement, Web, Mobile Devices

- Currently, >50,000 participants

- www.genesforgood.org

# Return of Results

# 10,000 Participants…



**Number of participants by gender and age group**

Number of participants 10000

| Age group | Male | Female |
|---|---|---|
| 71+ | 70 | 139 |
| 61-70 | 247 | 569 |
| 51-60 | 316 | 906 |
| 41-50 | 382 | 1 060 |
| 31-40 | 634 | 1 471 |
| 21-30 | 967 | 2 098 |
| under 21 | 186 | 415 |

Number of participants

■ Male  ■ Female

# BMI, Age & Diabetes



**Relationship of BMI with Diabetes Type 1 or 2**

Bays et al. (2007) *International Journal of Clinical Practice*

# Results: BMI GWAS

| Pheno | n | Chr:Pos | SNP | Gene | Our P | Other P* |
|-------|-----|---------|-----|------|-------|----------|
| BMI | 2,851 | 16:53803574 | rs1558902 | *FTO* | $5 \times 10^{-5}$ | $5 \times 10^{-120}$ |



*Speliotes et al. (2010) *Nature Genetics*

# Average Reported Sleep Hours Over a Year
## (data from Genes for Good participants)

Anita Pandit



2016

# Michigan Genomics Initiative

- Combine genetic and electronic health information on 40,000+ patients          50 new participants per day

- Use genetic information study many traits and diseases          Diverse traits – 40% w/cancer

- Build catalog of naturally occurring human knockouts          Speed and improve translation

- Clear, easy to understand consent – full participant buy-in.          Key for long term success

- **Team effort: Abecasis (Genetics), Ketherpal (Electronic Health Records), Brummett (Recruitment)**

# Michigan Genomics Initiative (Freeze 1)
## 20,000 individuals
## 7.5 million variants x 1,500 phenotypes

# Michigan Genomics Initiative Association Statistics
## http://pheweb.sph.umich.edu



near TCF7L2

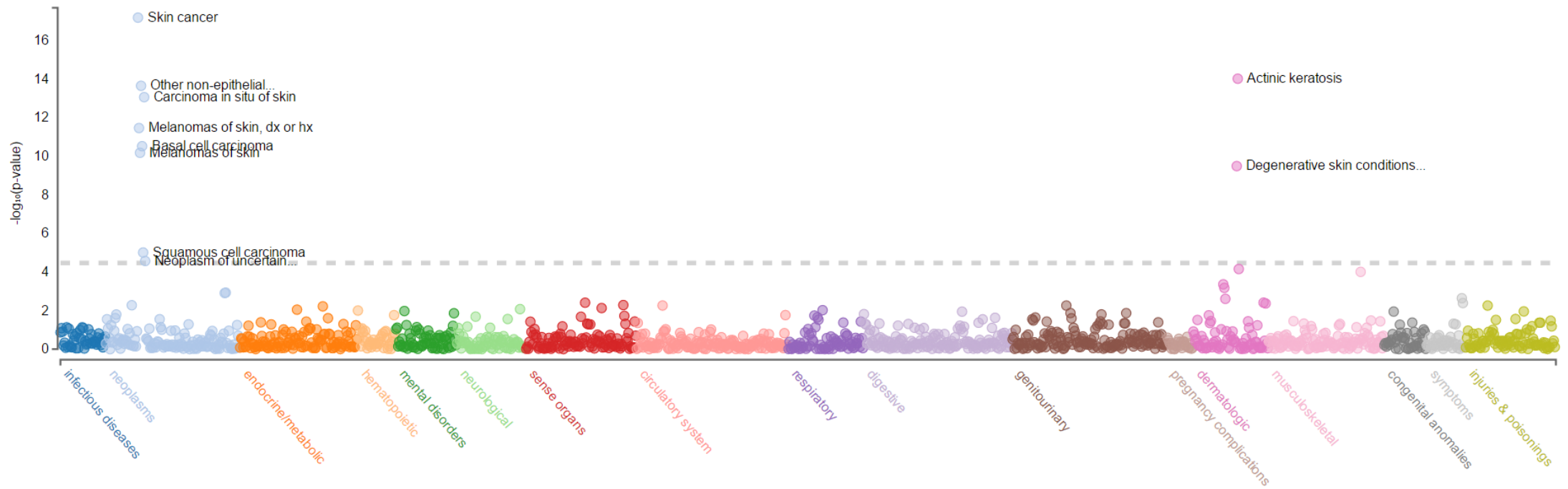Near HLA-DBQ1

# I heard rs10490924 in ARMS2 is associated with macular degeneration …
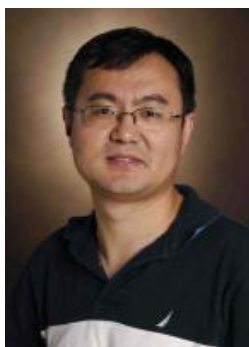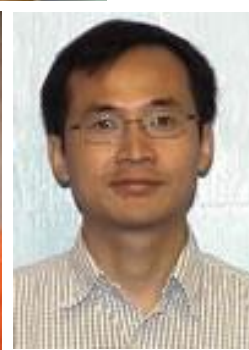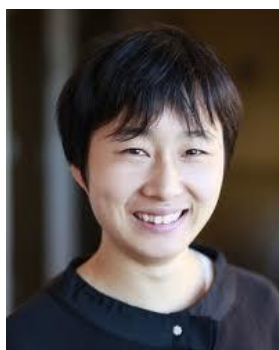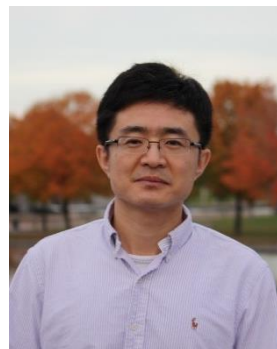
# I heard rs738409 in PNPLA3 is associated with liver disease ...

# I heard rs12203592 in IRF4 is associated with freckling, skin color ...

# The secret of success …

# Lessons learned…

- One person and a good idea can make a difference.

- The best students, postdocs, collaborators know something you don't.

- Take the time to be amazed. Drop everything and explore a new idea.

- Keep learning. There a so many great ideas out there.

- The most valuable tools and algorithms are often extremely simple.