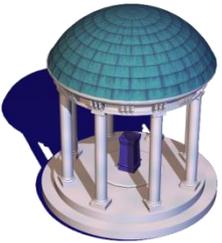




UNC
SCHOOL OF
PUBLIC HEALTH

FFGWAS

Fast **F**unctional **G**enome **W**ide **A**ssociation **A**nalysis
S of Surface-based Imaging Genetic Data



Chao Huang

Department of Biostatistics

Biomedical Research Imaging Center

The University of North Carolina at Chapel Hill,
Chapel Hill, NC 27599, USA

Joint work with Hongtu Zhu



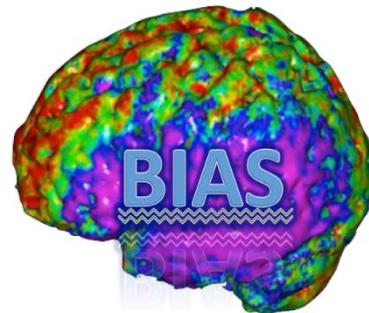


Outline

- **Motivation for Imaging Genetics**
- **Statistical Methods for Imaging Genetics**
- **Fast Functional Genome-wide Association Analysis**
- **Conclusion**

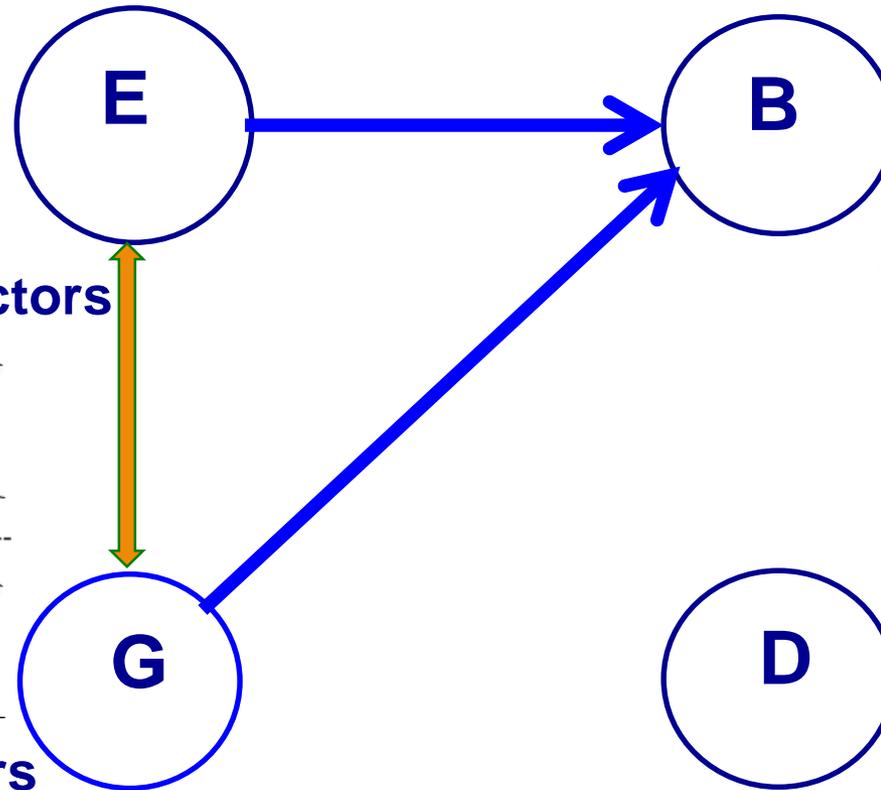
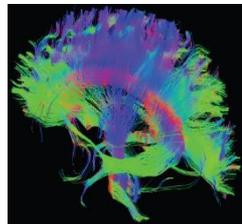
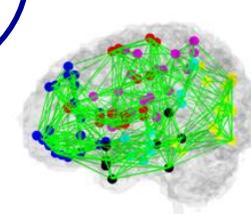
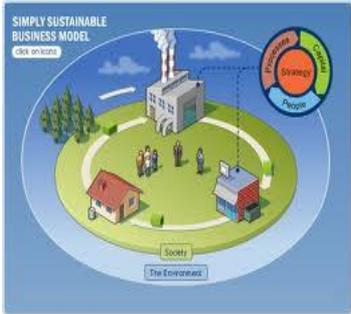


Motivation for Imaging Genetics

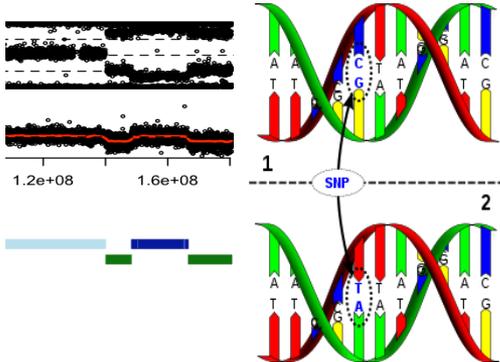




Imaging Genetics



E: environmental factors



G: genetic markers

D: disease

http://en.wikipedia.org/wiki/DNA_sequence



Imaging Data

**Structural
MRI**

- Variety of acquisitions
- Measurement basics
- Limitations & artefacts
- Analysis principles
- Acquisition tips

**Functional
MRI (task)**

**Diffusion
MRI**

**Functional
MRI
(resting)**

PET

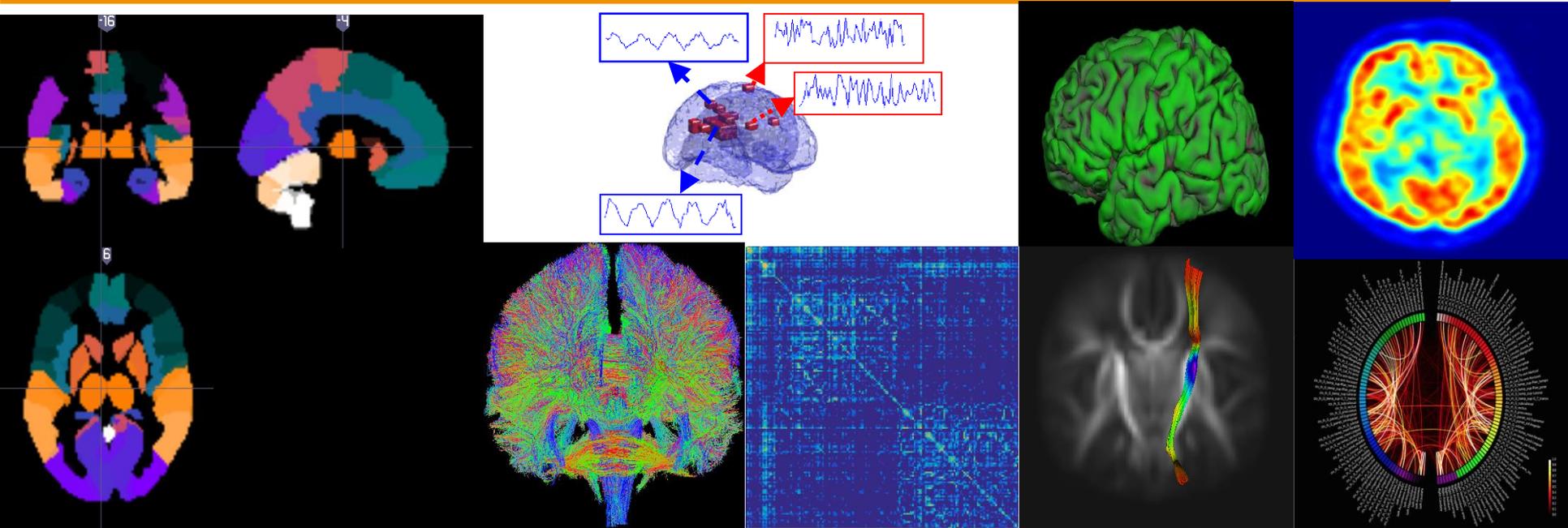
EEG/MEG

CT

Calcium



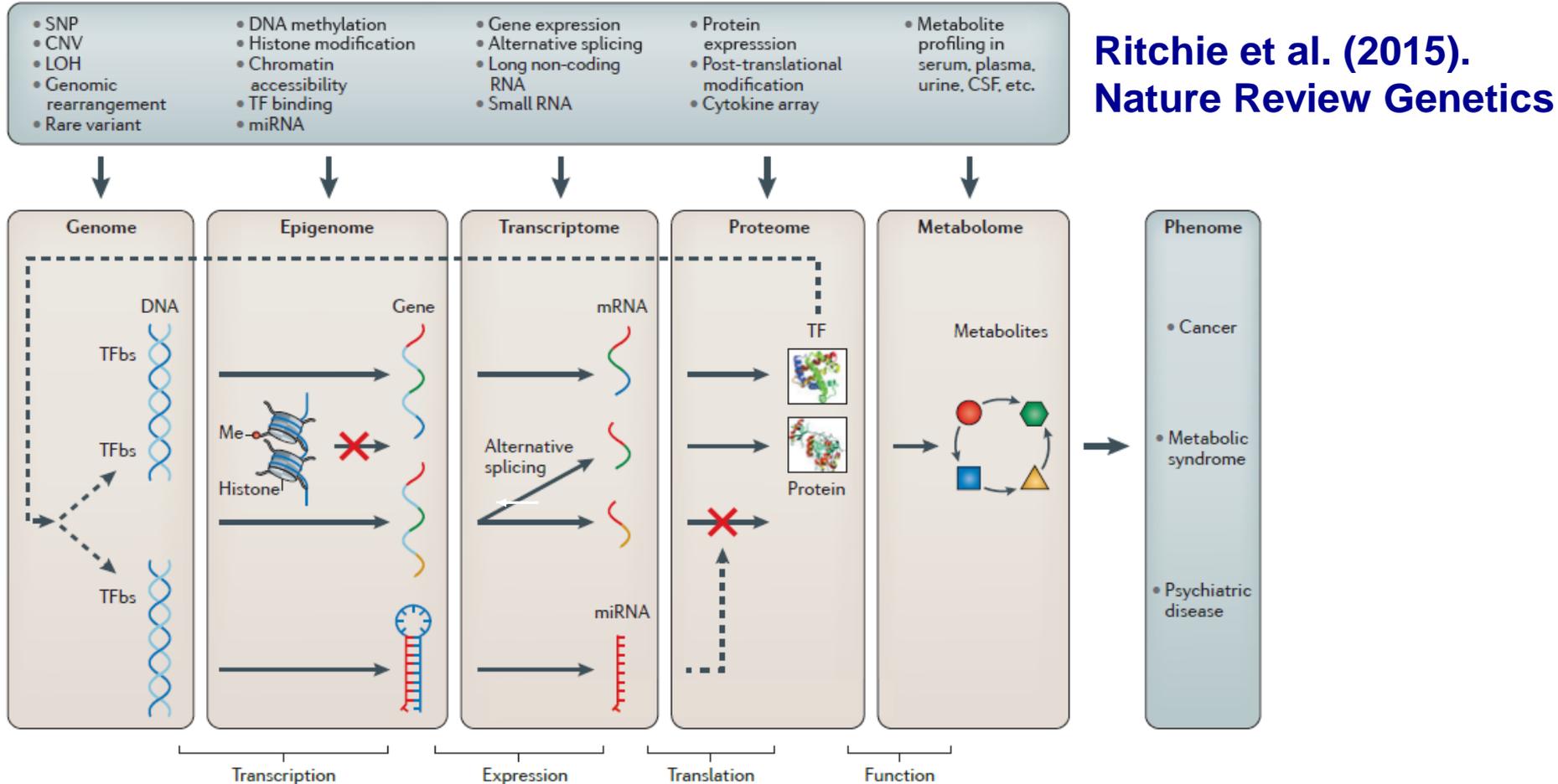
Neuroimaging Phenotype



Multivariate, smoothed functions, and piecewisely smoothed functions
Dimension varies from 100~500,000.



Multi-Omic Data



Ritchie et al. (2015).
Nature Review Genetics

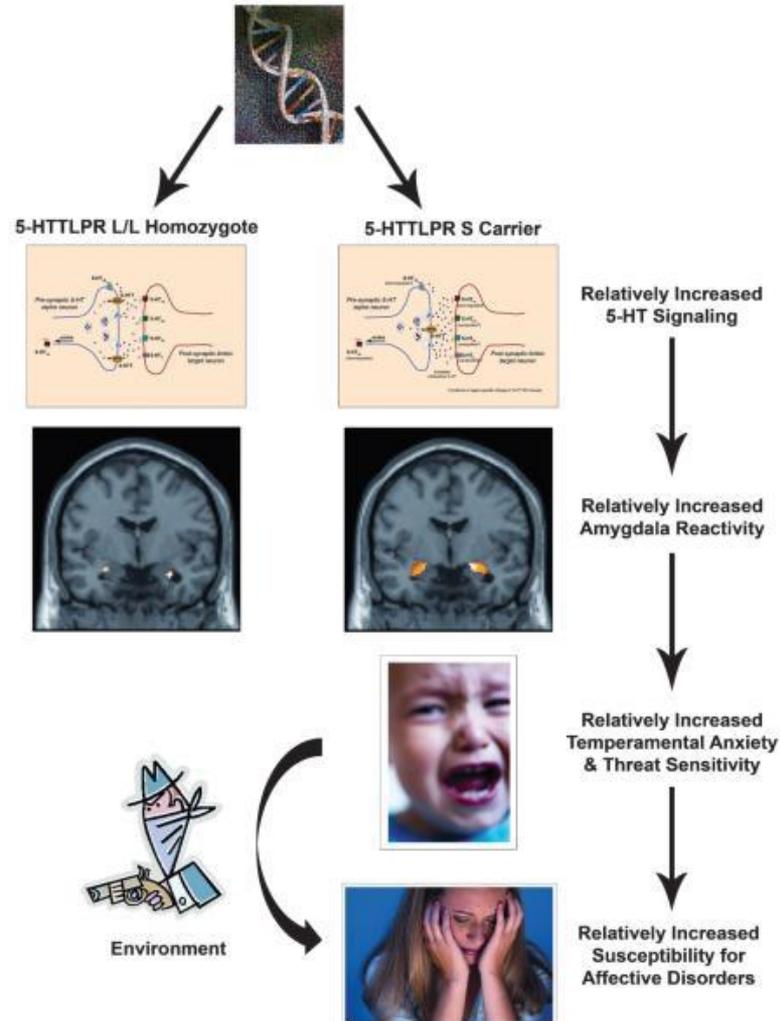


Motivation

Imaging genetics allows for the identification of **how common/rare genetic polymorphisms** influencing molecular processes (e.g., serotonin signaling), **bias neural pathways** (e.g., amygdala reactivity), **mediating individual differences in complex behavioral processes** (e.g., trait anxiety) related to disease risk in response to environmental adversity.

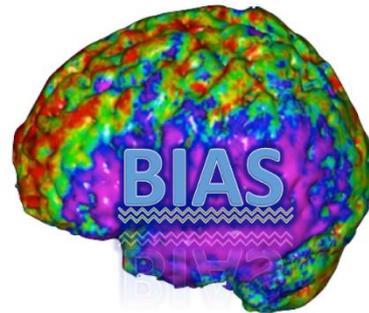
(Hariri AR, Holmes A.
Genetics of emotional regulation:
the role of the serotonin transporter in neural
function.

Trends Cogn Sci. [10:182–191])



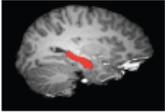
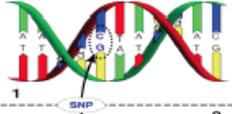


Statistical Methods for Imaging Genetics





Statistical Methods

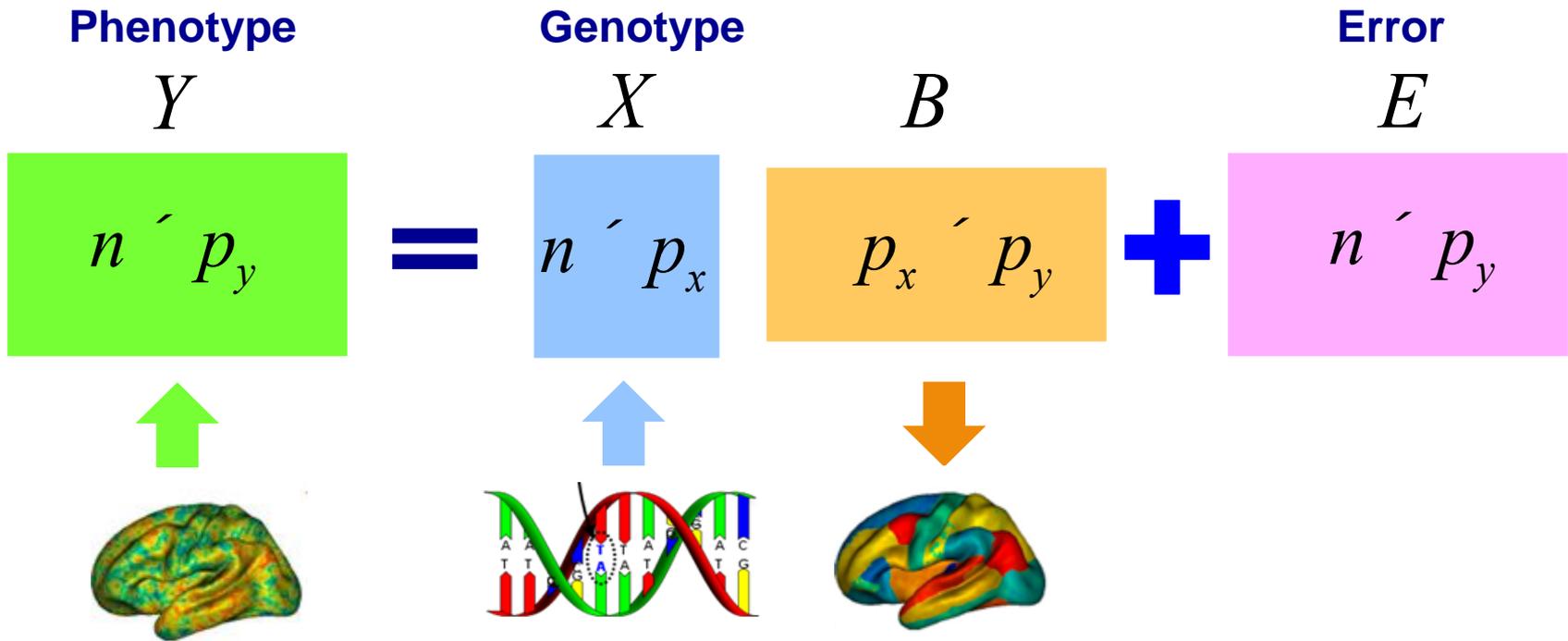
Imaging \ Genetics	Candidate ROI 	Many ROI 	Voxelwise 
Candidate SNP 	Imager	Imager	Imager
Candidate Gene 	Geneticist	↑	↑
Genome-wide SNP <pre>rs661983 rs59206197 r rs11493920 rs58524100 r rs34984204 rs11218322 r rs55682479 rs12279197 r rs664236 rs59966742 r rs34898405 rs617847 r</pre>	Geneticist	↑	↑
Genome-wide Gene <pre>BUD13 SCN4B CBL G BUD13 SCN2B MCAM G BUD13 AMICA1 MCAM G ZNF259 AMICA1 MFRP G ZNF259 AMICA1 MFRP G</pre>	Geneticist		



High Dimensional Regression Model

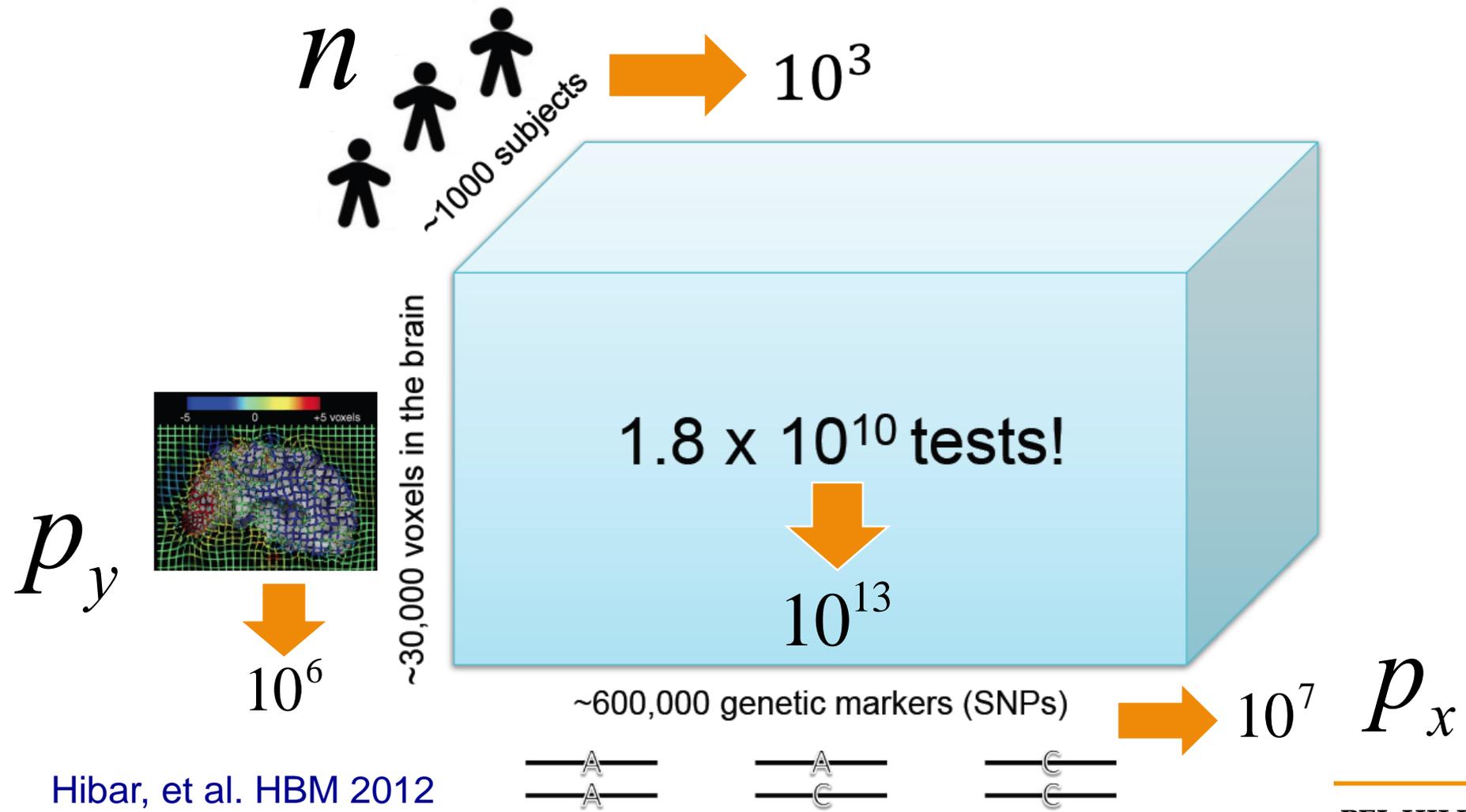
Data $\{(Y_i, X_i): i = 1, \dots, n\}$

$$Y_i = \{y_i(v): v \hat{=} V\} \quad X_i = \{X_i(g): g \hat{=} G_0\}$$

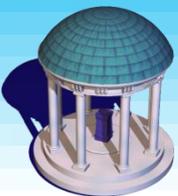




Challenges



Hibar, et al. HBM 2012



Fast Voxel-wise Genome-wide Analysis



Huang, et al. Neuroimage 2015

Issues to be addressed:

- Spatially correlated functional data
- Multivariate imaging phenotypes

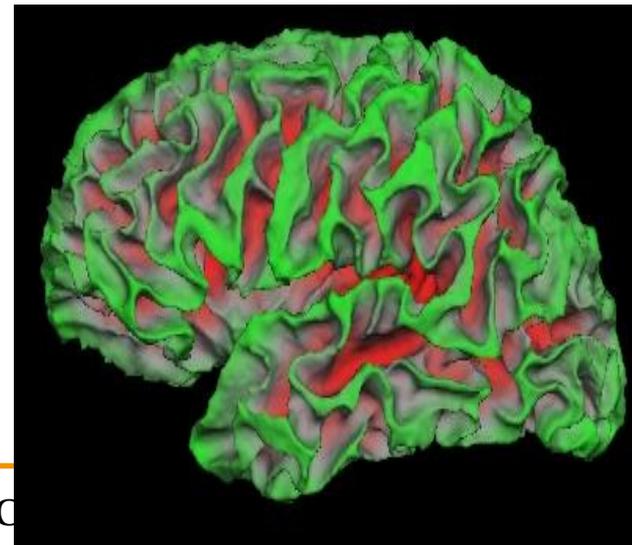
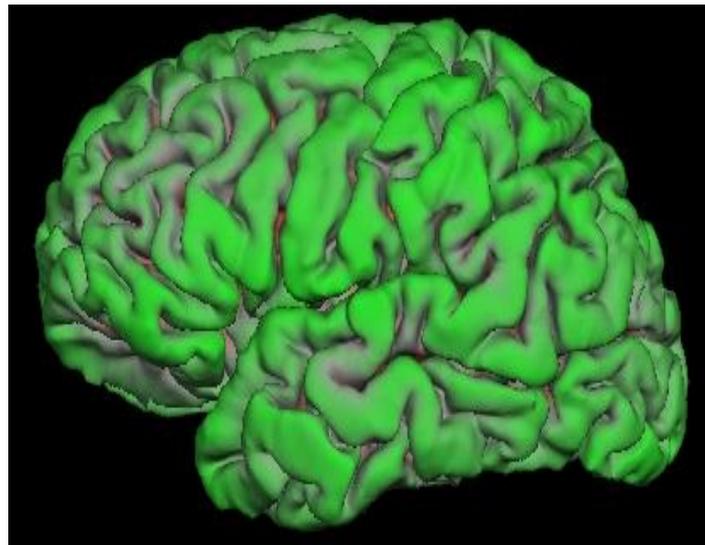
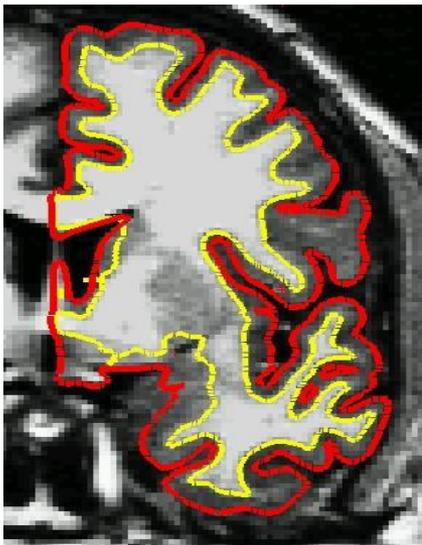
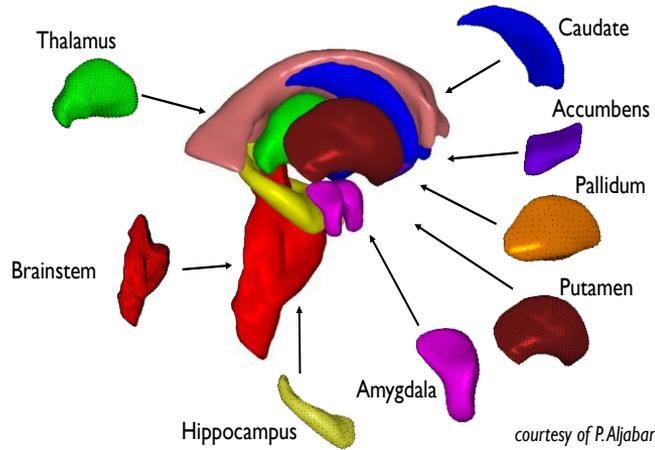
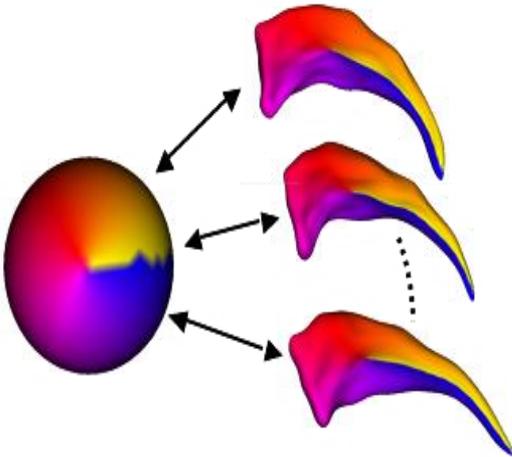


Fast Functional Genome-wide Analysis



Data Structure

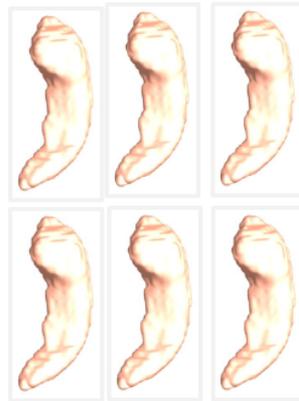
Have 15 different sub-cortical structures (left/right separately)



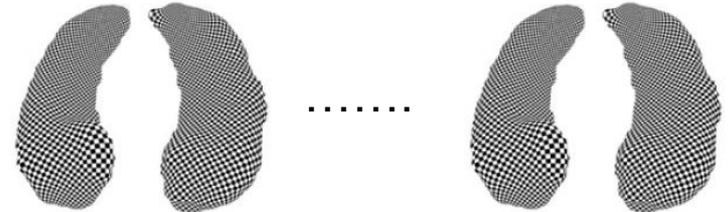


Data Structure

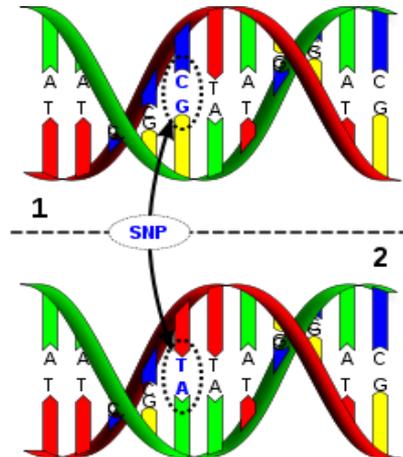
Hippocampal Surface



Person No.1 Person No.100



Genetic Variation



Person No. 1

Person No. 100

SNP1 SNP2 SNP2000

$$\begin{bmatrix} 1 & 2 & \cdots & 0 \\ 0 & \ddots & & 1 \\ \vdots & & \ddots & \vdots \\ 1 & 0 & \cdots & 2 \end{bmatrix}$$



FFGWAS

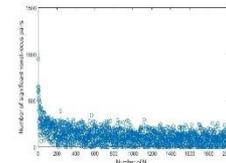
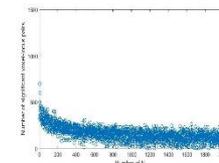
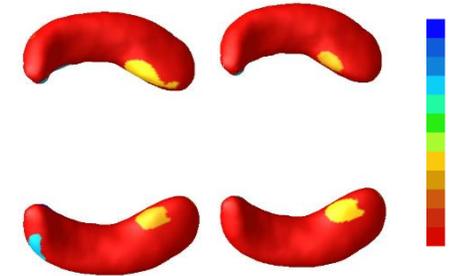
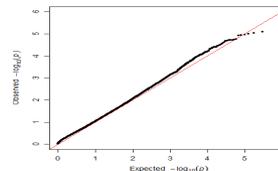
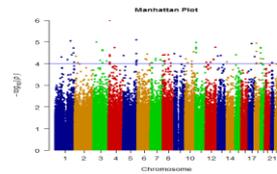
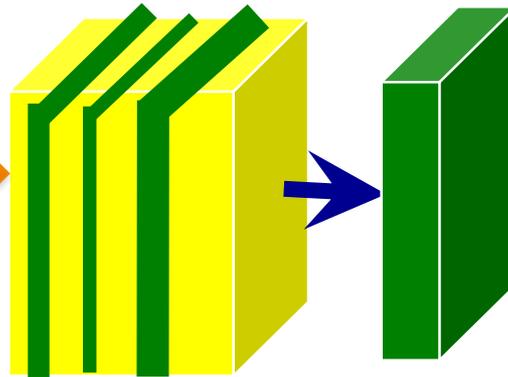
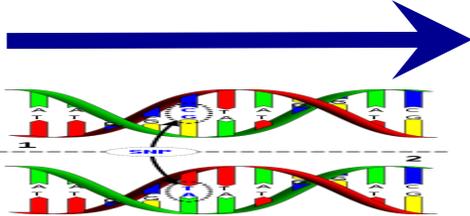
Multivariate
Varying
Coefficient
Model

Global Sure
Independence
Screening
Procedure

Significant
Voxel-locus
Cluster-locus
Detection



SNP 1 SNP N

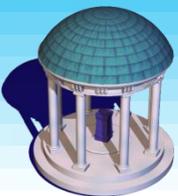




Multivariate Varying Coefficient Model

$$y_{i,j}(\mathbf{d}) = \mathbf{x}_i^T \boldsymbol{\beta}_j^{(c)}(\mathbf{d}) + \mathbf{z}_i(\mathbf{g})^T \boldsymbol{\beta}_j^{(g)}(\mathbf{d}) + \eta_{i,j}(\mathbf{d}) + \epsilon_{i,j}(\mathbf{d}), i = 1, \dots, n, j=1, \dots, J$$

where $\boldsymbol{\beta}_j^{(c)}(\mathbf{d})$ is a $p_x \times 1$ vector associated with non-genetic predictors (e.g., age, gender), and $\boldsymbol{\beta}_j^{(g)}(\mathbf{d})$ is an $p_g \times 1$ vector of genetic fixed effects (e.g., additive or dominant). Moreover, $\boldsymbol{\epsilon}_i(\mathbf{d}) = (\epsilon_{i,1}(\mathbf{d}), \dots, \epsilon_{i,J}(\mathbf{d}))^T$ are measurement errors and independent and identical copies of a stochastic process $\mathbf{SP}(\mathbf{0}, \boldsymbol{\Sigma}_\epsilon)$ and $\boldsymbol{\eta}_i(\mathbf{d}) = (\eta_{i,1}(\mathbf{d}), \dots, \eta_{i,J}(\mathbf{d}))^T$ are independent with $\boldsymbol{\epsilon}_i(\mathbf{d})$ and identical copies of a stochastic process $\mathbf{SP}(\mathbf{0}, \boldsymbol{\Sigma}_\eta^{(g)})$.



Multivariate Varying Coefficient Model

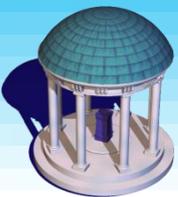
We need to test:

$$H_0 : \beta^{(g)}(\mathbf{d}) = 0 \text{ for all } \mathbf{d} \text{ v.s. } H_1 : \beta^{(g)}(\mathbf{d}) \neq 0$$

We first consider a local Wald-type statistic as:

$$T_n(g, \mathbf{d}) = \mathbf{r}^{(g)}(\mathbf{d})^T \left[\hat{\Sigma}_\eta^{(g)-1}(\mathbf{d}, \mathbf{d}) \otimes \left[\sum_{i=1}^n z_i(g)^{\otimes 2} \right]^{-1} \right] \mathbf{r}^{(g)}(\mathbf{d}),$$

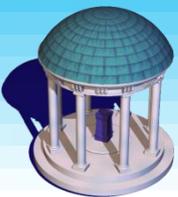
where $\mathbf{r}^{(g)}(\mathbf{d}) = \hat{\beta}^{(g)}(\mathbf{d}) - \text{Bias}(\hat{\beta}^{(g)}(\mathbf{d}))$.



Big-data Challenges

Several big-data challenges arise from the calculation of $T_n(g, d)$ as follows.

- Calculating $\hat{\Sigma}_\eta^{(g)}(d)$ across all loci and vertices can be computationally.
- Bandwidth selection in $T_n(g, d)$ across all loci can be also computationally.
- Holding all $T_n(g, d)$ in the computer hard drive requires substantial computer resources.
- Speeding up the calculation of $T_n(g, d)$.



FFGWAS

To solve these computational bottlenecks, we propose three solutions as follows.

- Calculate $\hat{\Sigma}_{\eta}^{(g)}(d)$ under the null hypothesis H_0 for all loci.
- Divide all loci into K groups based on their minor allele frequency (MAF), and select a common optimal bandwidth for each group.
- Develop a GSIS procedure to eliminate many ‘noisy’ loci based on a global Wald-type statistic.



A Global Sure Independence Screening

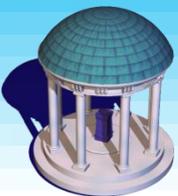
(1) The global Wald-type statistic at locus g is defined as

$$T_n(g) = \frac{1}{M} \text{tr} \left\{ \text{Vec}(\widetilde{X})^{\otimes 2} \left[\sum_{m=1}^M Y_w(d_m) \hat{\Sigma}_\eta^{-1}(d_m) Y_w^T(d_m) \right] \otimes \left[(\widetilde{X} \widetilde{X}^T)^{-1} [0_{p_g \times p_c} \quad I_{p_g}]^T \left[\sum_{i=1}^n z_i(g)^{\otimes 2} \right]^{-1} [0_{p_g \times p_c} \quad I_{p_g}] (\widetilde{X} \widetilde{X}^T)^{-1} \right] \right\}$$

(2) Calculate the p-values of $T_n(g)$ for all loci

(3) Sort the $-\log_{10}(p)$ -values of $T_n(g)$ and select the top N_0 loci

The candidate significant locus set $\tilde{\mathcal{G}}_0 = \{\tilde{g}_1, \dots, \tilde{g}_{N_0}\}$



Detection Procedure

(1) The first one is to detect significant voxel-locus pairs

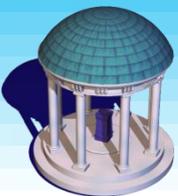
$$T_n(\tilde{g}, d) = \text{tr} \left\{ \text{Vec}(\tilde{X})^{\otimes 2} \left[Y_w(d) \hat{\Sigma}_\eta^{-1}(d) Y_w^T(d) \right] \otimes \left[(\tilde{X} \tilde{X}^T)^{-1} [0_{p_g \times p_c} \ I_{p_g}]^T \left[\sum_{i=1}^n z_i(g)^{\otimes 2} \right]^{-1} [0_{p_g \times p_c} \ I_{p_g}] (\tilde{X} \tilde{X}^T)^{-1} \right] \right\}$$

(2) The second one is to detect significant cluster-locus pairs.

Wild Bootstrap method



Simulation Studies and Real Data Analysis



Simulation Studies: Data Generation

Covariate Data (non-genetic data)

Generated from either $U(0,1)$ or the Bernoulli distribution with success probability 0.5.

Genetic Data

Linkage Disequilibrium (LD) blocks (**Haploview** & **PLINK**)

1. Generate 2,000 blocks;
2. Randomly select 10 SNPs in each block;
3. Chose the first 100 SNPs as the causal SNPs



Simulation Studies: Data Generation

Imaging Data

Step 1: Fitting the model without genetic predictors

$$y_{i,j}(d) = x_i^T \beta_j^{(c)}(d) + z_i^{(g)T} \beta_j^{(g)}(d) + \eta_{i,j}(d) + \epsilon_{i,j}(d), i = 1, \dots, n, j = 1, \dots, J$$

→ Estimates of $\beta_j^{(c)}(d)$ Σ_ϵ $\Sigma_\eta^{(g)}$ → True values

Step 2: Specifying effected Regions Of Interest associated with causal SNPs

$$\beta_j^{(g)}(d) = \begin{cases} r, & \forall d \in \text{ROIs} \\ 0, & \text{otherwise} \end{cases}$$

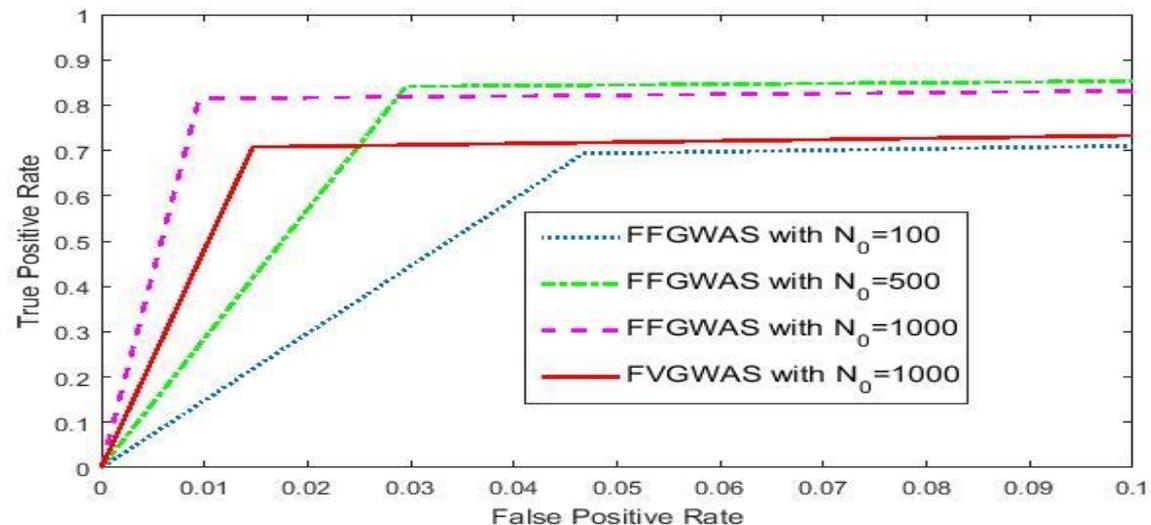
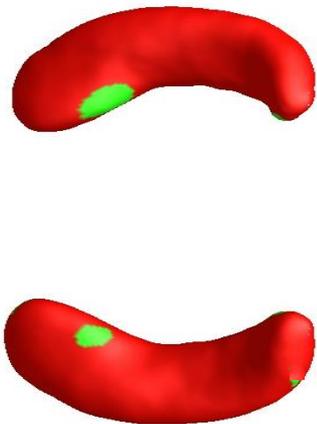
Step 3: Generating imaging data with prespecified parameters and ROIs



Simulation Studies

Simulation settings: the green and red regions in the figure, respectively, represent Hippocampal surface, and the effected ROI associated with the causal SNPs among first 20000 SNPs.

$$\beta_j^{(g)}(d) = 0.001$$



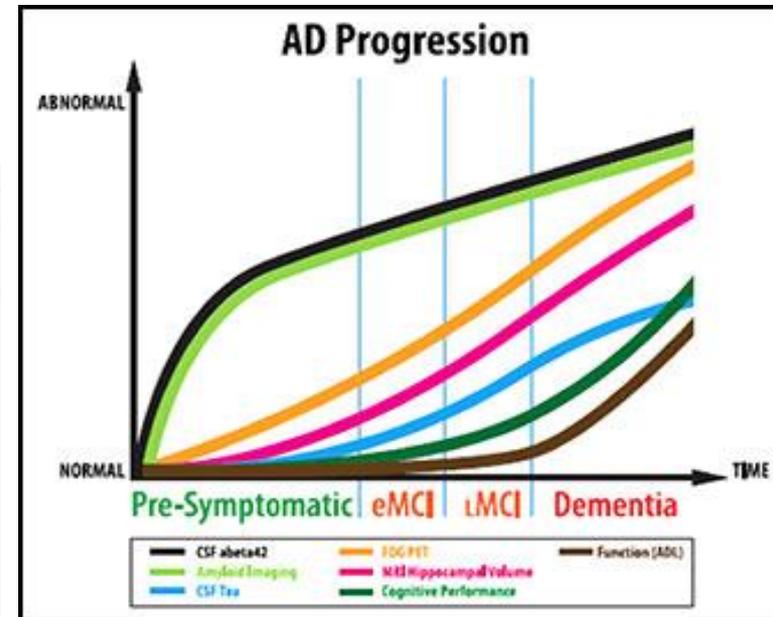
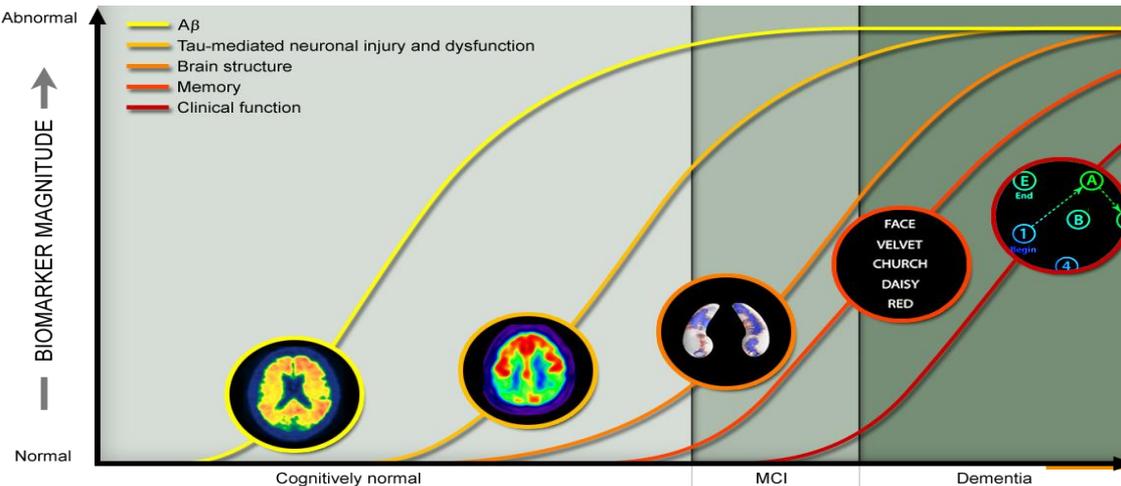
Simulation results for comparisons between FFGWAS and FVGWAS in identifying significant voxel-SNP pairs.

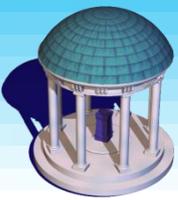


Imaging Genetics for ADNI

PI: Dr. Michael W. Weiner

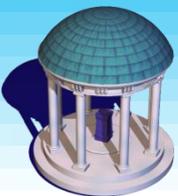
- detecting AD at the earliest stage and marking its progress through biomarkers;
- developing new diagnostic methods for AD intervention, prevention, and treatment.
- A longitudinal prospective study with 1700 aged between 55 to 90 years
- Clinical Data including Clinical and Cognitive Assessments
- Genetic Data including Illumina SNP genotyping and WGS
- MRI (fMRI, DTI, T1, T2)
- PET (PIB, Florbetapir PET and FDG-PET)
- Chemical Biomarker





ADNI Data Analysis: Dataset Description

- **708** MRI scans of AD (**186**), MCI (**388**), and healthy controls (**224**) from ADNI-1.
- These scans on **462** males and **336** females are performed on a 1.5 T MRI scanners.
- The typical protocol includes the following parameters:
 - (i) repetition time (TR) = 2400 ms;
 - (ii) inversion time (TI) = 1000 ms;
 - (iii) flip angle = 8° ;
 - (iv) field of view (FOV) = 24 cm with a 256 x 256 x 170 acquisition matrix in the x-, y-, and z-dimensions,
 - (v) voxel size: 1.25 x 1.26 x 1.2 mm³.
- Covariates: gender, age, APOE $\epsilon 4$, and the top 5 PC scores in SNPs



Imaging Data Preprocessing

Surface fluid registration based hippocampal sub-regional analysis package (Shi et al., Neuroimage, 2013)

- **Hippocampal surface registration**
isothermal coordinates and fluid registration
- **Surface statistics computation**
 1. multivariate tensor-based morphometry (mTBM) statistics
 2. radial distance

Finally, we obtained left and right hippocampus shape representations as 100×150 matrices.



ADNI Data Analysis

Top 10 SNPs (Left Hippocampus)

Top 10 SNPs (Right Hippocampus)

SNP	CHR	BP	-LOG 10(p)
rs657132	18	2.20533e+07	7.579767
rs604345	18	2.20033e+07	6.729377
rs582110	18	2.19954e+07	6.672876
rs546000	18	2.20031e+07	6.672876
rs489631	18	2.1989e+07	6.620395
rs16837577	1	1.94871e+08	6.016773
rs3812872	13	6.19869e+07	5.468391
rs6826085	4	7.68702e+07	5.459163
rs929714	7	1.3263e+08	5.314317
rs2042067	7	1.32651e+08	5.306583

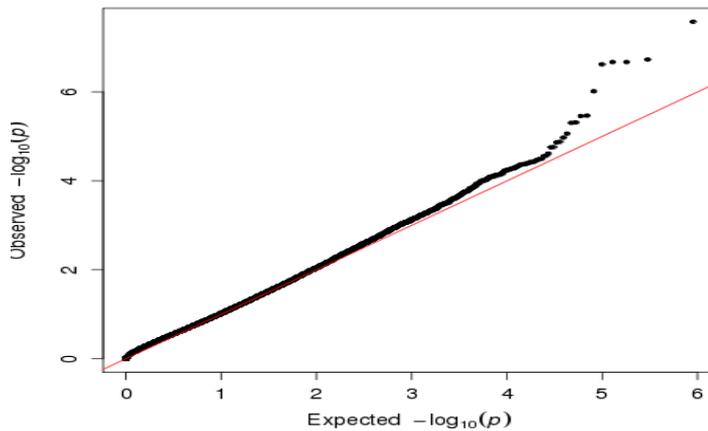
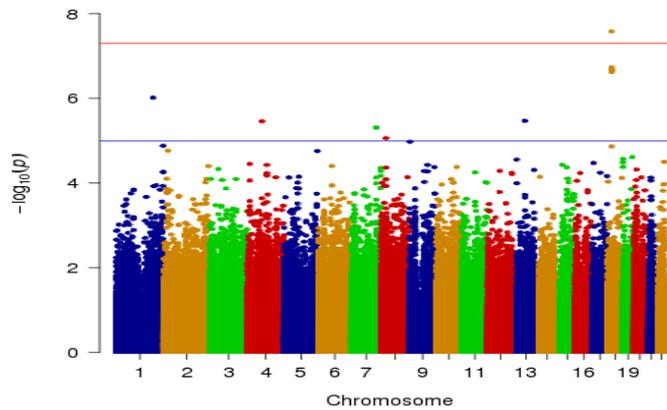
SNP	CHR	BP	-LOG 10(p)
rs4681527	3	1.44e+08	6.764886
rs3108514	2	1.51279e+08	6.274511
rs12264728	10	1.3214e+08	5.961976
rs652911	10	1.3214e+08	5.739661
rs10801705	1	8.95004e+07	5.622668
rs366346	10	1.32141e+08	5.617185
rs7312068	12	2.94352e+07	5.604041
rs7617465	3	1.43999e+08	5.522112
rs17605251	7	1.02746e+08	5.486603
rs749788	2	2.84618e+06	5.474675



ADNI Data Analysis

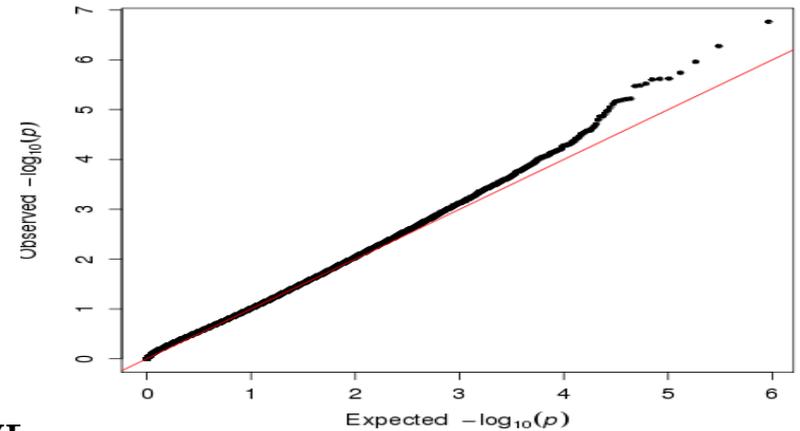
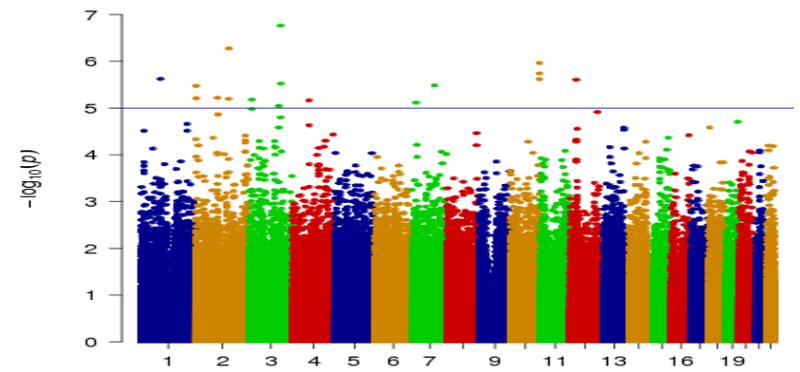
(Left Hippocampus)

Manhattan Plot



(Right Hippocampus)

Manhattan Plot

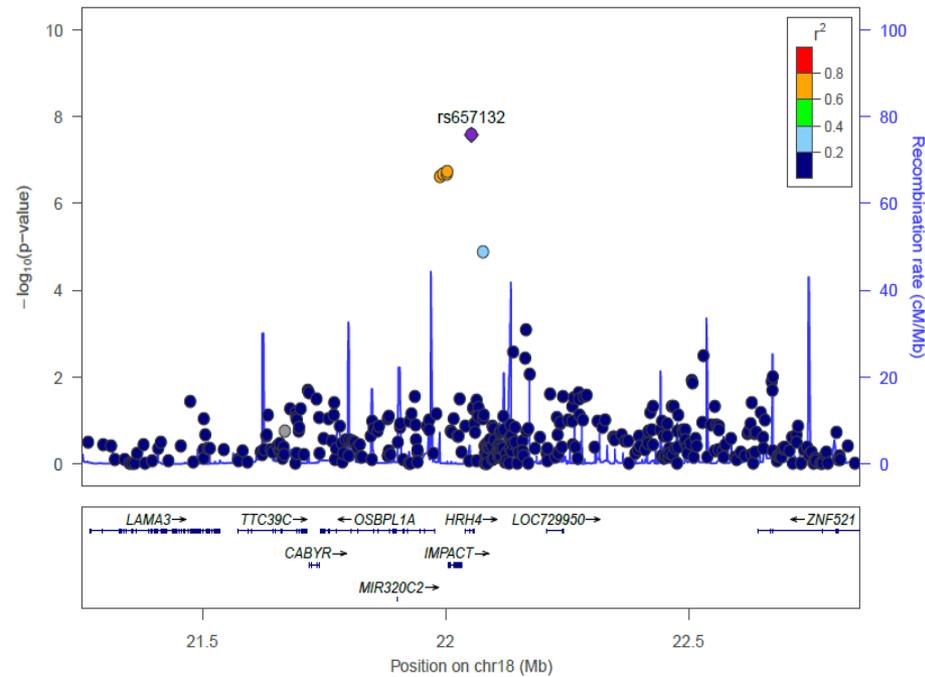




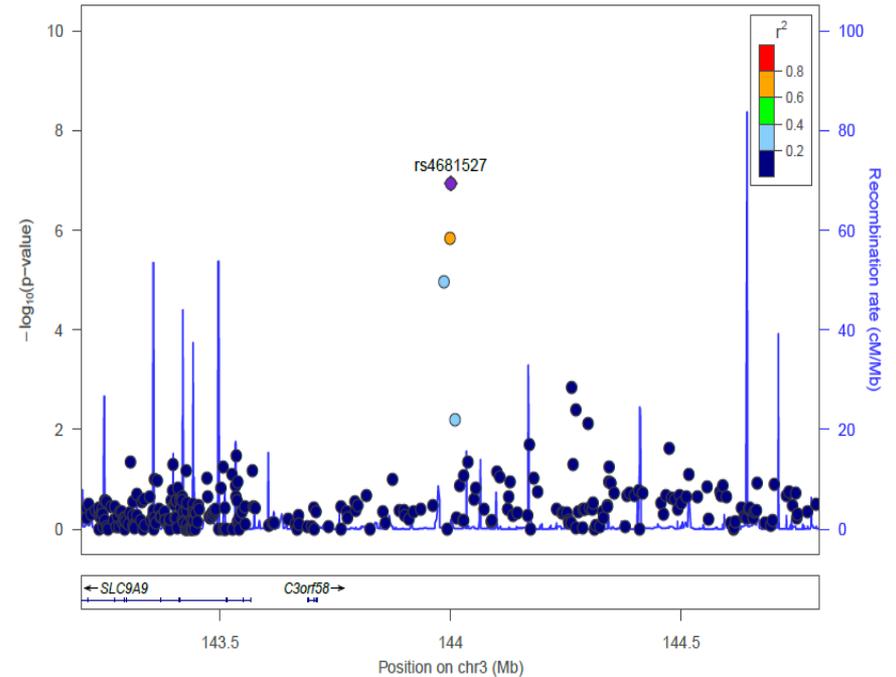
ADNI Data Analysis: Left Hippocampus

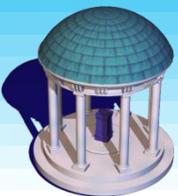
Significant Loci Zoom

(Left Hippocampus)



(Right Hippocampus)





ADNI Data Analysis: Left Hippocampus

(Left Hippocampus)

Top 1 SNP: rs657132

Closed Gene: HRH4

HRH4 (Histamine Receptor H4) is a Protein Coding gene.

Diseases associated with HRH4: cerebellar degeneration

An important paralog of this gene: CHRM4

Mirshafiey & Naddafi, Am J Alzheimers Dis Other Demen. 2013

(Right Hippocampus)

Top 1 SNP: rs4681527

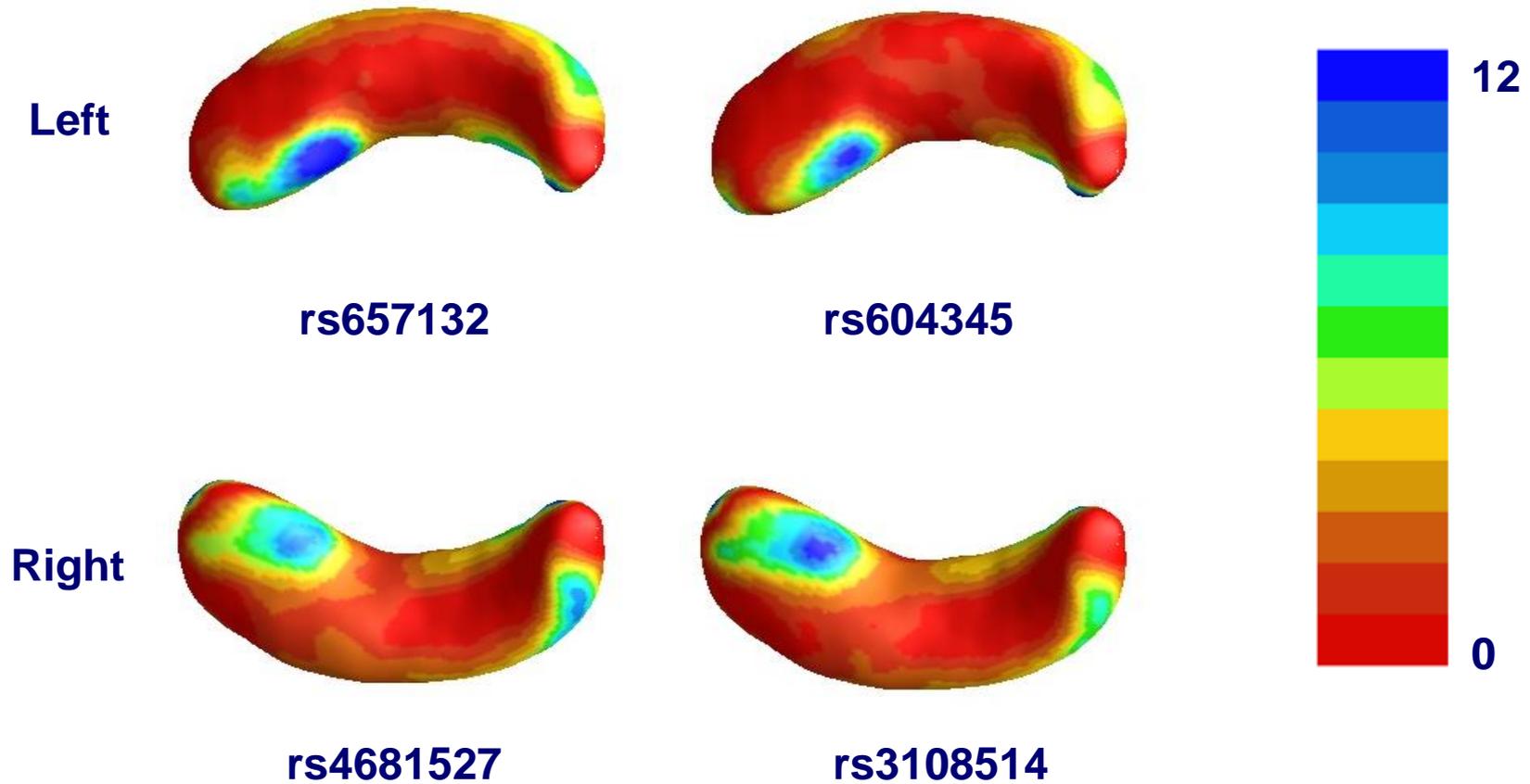
Closed Gene: C3orf58

C3orf58 (Chromosome 3 Open Reading Frame 58) is a Protein Coding gene.

Diseases associated with C3orf58: hypoxia



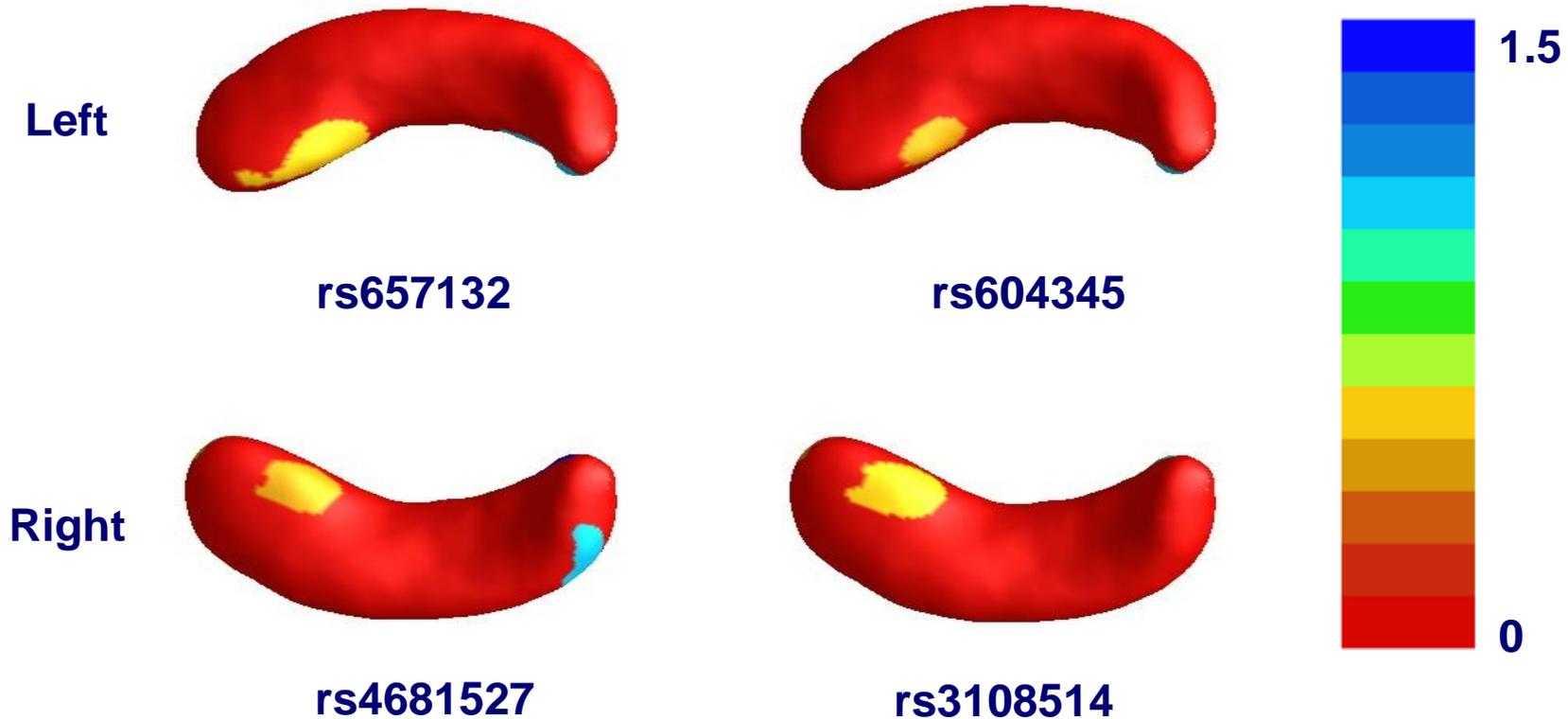
ADNI Data Analysis



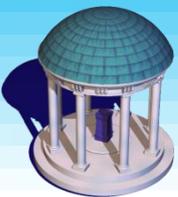
$-\log_{10}(p)$ values on Hippocampus (L & R) corresponding to Top 2 SNPs



ADNI Data Analysis



$-\log_{10}(p)$ values of significant clusters on Hippocampus (L & R) corresponding to Top 2 SNPs



Conclusion

- **We have developed a FFGWAS pipeline for efficiently carrying out whole-genome analyses of multimodal imaging data.**
- **Our FFGWAS consists of a multivariate varying coefficient model, a global sure independence screening (GSIS) procedure, and a detection procedure based on wild bootstrap methods.**
- **Two key advantages of using FFGWAS include**
 - (i) Much smaller computational complexity;
 - (ii) GSIS for screening many noisy SNPs.
- **We have successfully applied FFGWAS to hippocampal surface data & genetic data of ADNI study.**



A Software for FFGWAS

