# CONTENTS
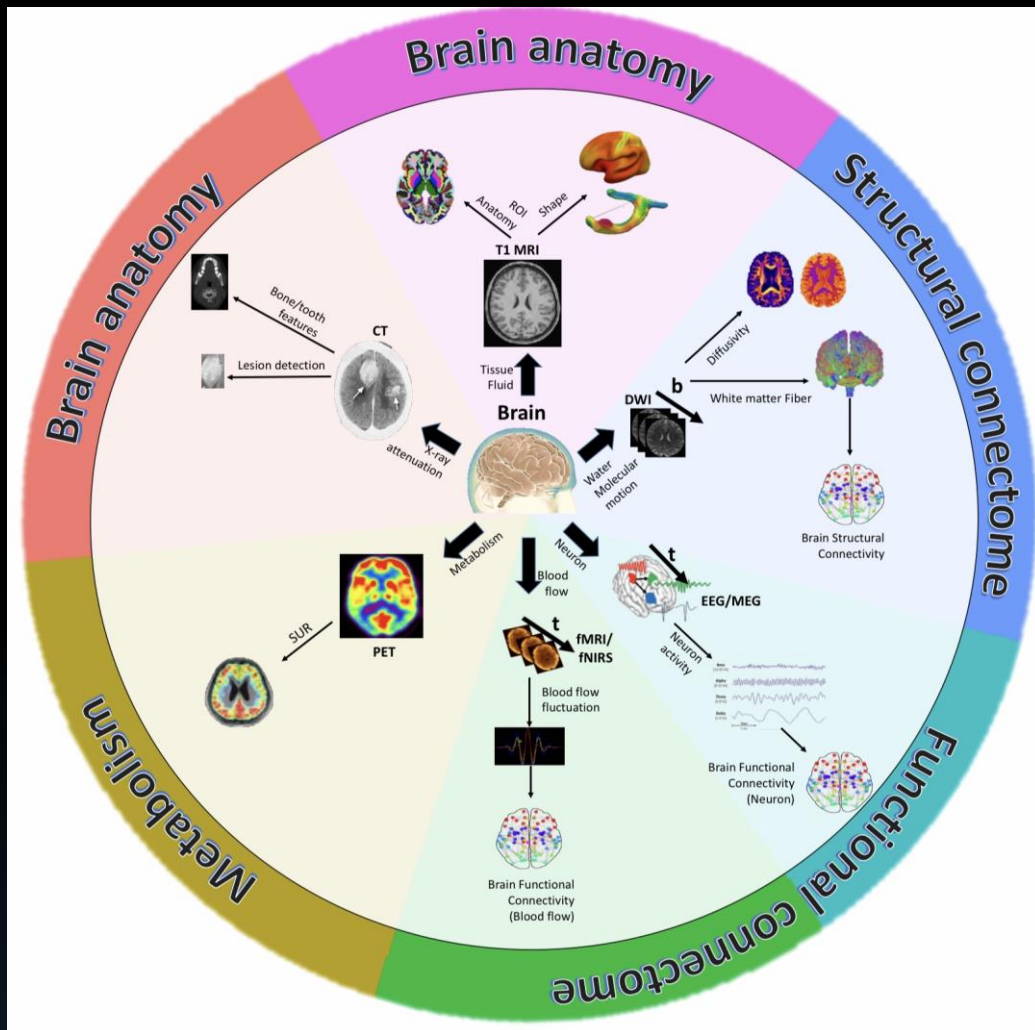
# Part I

## A Review of Neuroimaging Techniques

# Eight Popular Neuroimaging Techniques



- Structural magnetic resonance imaging (sMRI)
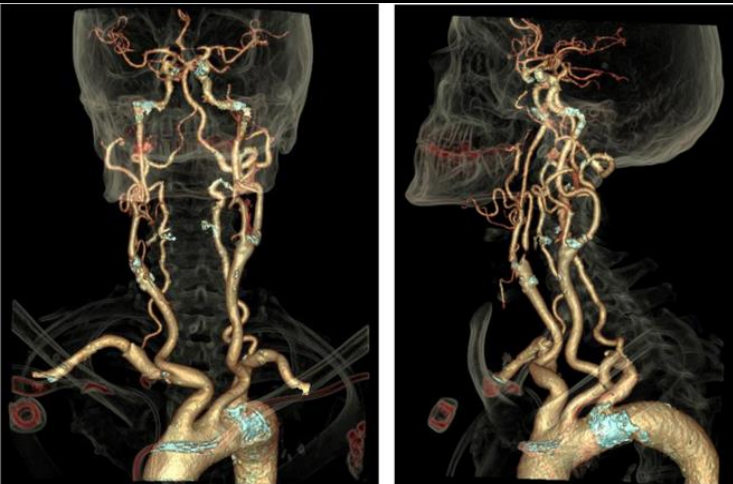- Diffusion weighted MRI (DWI)
- Functional MRI (fMRI)
- Positron emission tomography (PET)
- Computerized tomography (CT)
- Electroencephalography (EEG)
- Magnetoencephalography (MEG)
- Functional near-infrared spectroscopy (fNIRS)

Each image modality has its tracer, data dimension, extracted features, and  main clinical and research applications.

# CT, PET, MEG, EEG, and fNIRS



https://www.omegapds.com/ct-angiography-of-the-head-or-neck/

https://en.wikipedia.org/

# sMRI, fMRI, and DWI

# A Multi-model Approach



Image by A. Galka





**The van Essen diagram**

- Different models at different scales.
- Ladder of overlapping models.
- Must be testable against multiple phenomena.

# Part II

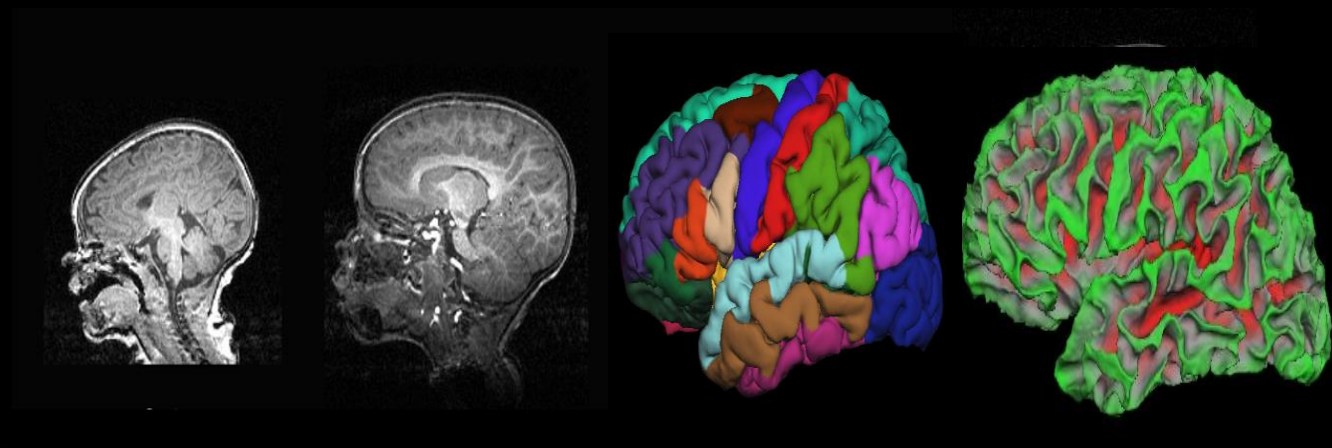## Imaging Processing Analysis Methods

# Four Common Themes (CT1)-(CT4)

**(CT1)** Complex Brain Objects



https://dana.org/article/neuroanatomy-the-basics/

**(CT2)** Complex Spatiotemporal Structures
- spatial and temporal resolutions
- spatio-temporal smoothness
- spatiotemporal correlation

**(CT3)** Extremely High Dimensionality



**(CT4)** Heterogeneity within Individual Subjects and across Centers/Studies

# Image Models

Image= $f$(B(age, gene, race, disease, others), device, acquisition, noises)

# Image Processing Analysis Methods

## IPA: Deconvolution
➤ **Image Reconstruction Process**
➤ **Image Enhancement Process**



## IPA: Structural Learning
➤ **Image Segmentation Process**
➤ **Image Registration Process**



Example: Airway Segmentation from CT

**Image Reconstruction Process (IRP)**



$$\frac{S(\mathbf{q})}{S_0} = \int P(\mathbf{r}, \Delta)e^{i\mathbf{q}\cdot\mathbf{r}}d\mathbf{r}; \quad \mathbf{q} = \gamma\delta g\mathbf{u}$$

## Image Enhancement Process (IEP)

- ❖ **Denoising**
- ❖ **Super-resolution**
- ❖ **Bias-field correction**
- ❖ **Harmonization**

# IPA: Structural Learning-ISP

**Image Segmentation Process (ISP)**

❖ **Quantification of brain development**

❖ **Localization of pathology**

❖ **Surgical planning**

❖ **Image-guided interventions**

❖ **Computer-aided detection and diagnosis**

❖ **Brain parcellation**



| Method | Dice_ET | Dice_WT | Dice_TC |
|---|---|---|---|
| Phase1 | 0.75245 | 0.89571 | 0.81561 |
| Phase2 Model1 | 0.75983 | 0.90397 | 0.82489 |
| Phase2 Model2 | 0.76091 | 0.90616 | 0.83622 |
| Phase2 Model3 | 0.74669 | 0.90349 | 0.8278 |
| Phase2 Model4 | 0.74187 | 0.90435 | 0.83211 |
| Phase2 Model5 | 0.75779 | 0.90733 | 0.83824 |
| Phase2 Model6 | 0.76091 | 0.90420 | 0.83713 |
| Phase2 Model7 | 0.76814 | 0.90574 | 0.84704 |
| Phase2 Model8 | 0.75440 | 0.90594 | 0.83826 |
| Phase2 Model9 | 0.78582 | 0.90491 | 0.83689 |
| XGBoost+ | 0.80536 | 0.91044 | 0.85057 |

# IPA: Structural Learning-IRP

**Image Registration Process (IRP)**

- **Automated image segmentation**
- **Construction of brain atlas**
- **Localization of pathology**
- **Multimodal fusion**
- **Population analysis**
- **Quantification of brain development**
- **Shape analysis**



**Fiber atlas**

**Fiber skeleton atlas**

Posterior-Multimodal
Cingulo-Opercular
Dorsal-Attention
Somatomotor
Language
Default
Visual1
Visual2
Auditory
Frontoparietal
Ventral-multimodal
Orbito-Affective

# Brain Function-based Structural Connectome Atlas



Stage 1:

Whole-brain Structural Connectome

Stage 2:

Creation of Fiber Skeleton

Stage 3:

Sparse Representation

# Longitudinal Elastic Shape Data Analysis



Figure 2: A schematic

# Challenges

**There is no publicly available, high-quality neuroimaging datasets with detailed annotation information that cover a large spectrum of segmentation tasks  in neuroimaging research**

# Part III

## Large-scale Neuroimaging related Studies

# Large-scale Neuroimaging-related Studies

PING - 900 Pediatric Imaging, Neurocognition, and Genetics
BCP - 300 Baby Connectome Project
ADNI - 2000 Alzheimer's Disease Neuroimaging Initiative
PNC - 1400 Philadelphia Neurodevelopmental Cohort
HCP - 1200 Human Connectome Project
ABCD - 10000 Adolescent Brain Cognitive Development
UKB - 500,000 UK Biobank Project
TCIA – 37,600 The Cancer Imaging Archive
NLST - 19,000 National Lung Screening Trial
OAI – 4800 Osteoarthritis Initiative

# Alzheimer's Disease Neuroimaging Initiative

The overall goal of ADNI is to validate potentially useful biomarkers for AD clinical treatment trials. ADNI is a multisite, prospective clinical study and actively supports the investigation and development of treatments that may slow or stop the progression of AD https://adni.loni.usc.edu/study-design. Researchers across 63 sites in the US and Canada have been tracking the progression of AD through clinical, imaging, genetic and biospecimen biomarkers, starting from normal aging, early mild cognitive impairment (EMCI), late mild cognitive impairment (LMCI) to dementia or AD.



**2004-now**

# The Human Connectome Project and Beyond

**The primary goals of HCP  include**
- **building a ``network map'' that will shed light on the anatomical and functional connectivity within the healthy human brain,**
- **promoting the understanding of inter-individual variability of brain circuits to behavior,**
- **facilitating research into brain disorders, such as  autism, AD, and schizophrenia, and**
- **making all data freely available to the scientific community.**

*The Heavily Connected Brain*
Peter Stern, "Connection, connection, connection…", Science, Nov. 1 2013: Vol. 342 no. 6158 P.577

**The NIH Human Connectome Project**
- **The Harvard/MGH-UCLA project**
- **The WU-Minn Project**

**The EU's 7th Framework Programme for Research**
- **Consortium Of Neuroimagers for the Non-Invasive Exploration of Brain Connectivity and Tracts**

**Healthy Adult Connectome**

**Lifespan Connectome Data**

**Connectomes related to Human Disease**

# Adolescent Brain Cognitive Development

The ABCD study is the largest prospective longitudinal study of brain development and child health in the United States, which has recruited approximately 11,880 children aged 9-10 years old from 21 research sites and is following them for 10 years into early adulthood.

Its initial goal was to examine risk and resiliency factors associated with the development of substance use, and then expanded far beyond, into identifying the underlying biospecimens, neural alterations, and environmental factors, and their contributions to the development of behavior, brain function, and other mental and physical outcomes throughout adolescence.



**2015-now**

https://abcdstudy.org/

# The UK Biobank Study

UK Biobank has collected and continues to collect extensive environmental, lifestyle, and genetic data on half a million participants.



**2006-now**



- **Imaging:** Brain, heart and full body MR imaging, plus full body DEXA scan of the bones and joints and an ultrasound of the carotid arteries. The goal is to image 100,000 participants, and to invite participants back for a repeat scan some years later.
- **Genetics:** Genotyping, whole exome sequencing & whole genome sequencing for all participants.
- **Health linkages:** Linkage to a wide range of electronic health-related records, including death, cancer, hospital admissions and primary care records.
- **Biomarkers:** Data on more than 30 key biochemistry markers from all participants, taken from samples collected at recruitment and the first repeat assessment.
- **Activity monitor:** Physical activity data over a 7-day period collected via a wrist-worn activity monitor for 100,000 participants plus a seasonal follow-up on a subset.
- **Online questionnaires:** Data on a range of exposures and health outcomes that are difficult to assess via routine health records, including diet, food preferences, work history, pain, cognitive function, digestive health and mental health.
- **Repeat baseline assessments:** A full baseline assessment is undertaken during the imaging assessment of 100,000 participants.
- **Samples:** Blood & urine was collected from all participants, and saliva for 100,000.

# ENIGMA

The major goals of ENIGMA include
- pushing forward the field of imaging genetics,
- ❖ ensuring promising and reproducible findings,
- ➢ sharing data, ideas, methods, algorithms and other information, and
- training new investigators.

The Enhancing NeuroImaging Genetics through Meta-Analysis (ENIGMA) Consortium is a global alliance of over 1,400 scientists across 43 countries in the fields of imaging genomics, neurology, and psychiatry, studying a range of large-scale human brain studies that integrate data based on sMRI, DWI, fMRI, genetic data and many patient populations from over 70 institutions worldwide

https://enigma.ini.usc.edu/

**2009-now**

Part IV

**Population-based Statistical Analysis Methods**

# Four Common Themes (CT5)-(CT8)

**(CT5) Sampling Bias**
- undercoverage
- observer bias,
- voluntary response bias
- survivorship bias
- recall bias
- exclusion bias

**(CT6) Complex Missing Data Patterns**
- ❖ missing by design
- ❖ faulty scanning
- ❖ attrition in longitudinal studies
- ❖ mis-entry
- ❖ non-responses in surveys

**(CT7) Complex Data Objects**



**(CT8) Complicated Causal Pathways in Brain-related Disorders**

# Population-based Statistical Analysis (PSA)

❖ **Study Design**
❖ **Statistical Parametric Mapping**
❖ **Object Oriented Data (OOD) Analysis**
❖ **Imputation Methods**
❖ **Data Integration Methods**



➢ **Dimension Reduction Methods**
➢ **Image Genetics**
➢ **Causality Research**
➢ **Predictive Analysis**
➢ **Knowledge-based Methods**

# Study Design

❖ **Case control study**

❖ **Cross-sectional study**

❖ **Cohort study**

❖ **Experimental study**

❖ **Descriptive Study:  case reports, case series, Descriptive surveys.**

➢ The UKB is a large, population-based cohort study, and many cross-sectional analyses have been conducted based on baseline data from UKB.

➢ The UKB is well known for its "healthy volunteer" selection bias, and may not be a true representation of the general population.

➢ Neuroimaging biomarkers are usually secondary outcome.

# Statistical Parametric Mapping

**Univariate Statistics**

Preprocessed data: single voxel

Design matrix

RFT/ permutation

➢ **More complex models**
➢ **Multiple comparisons**

Parameter estimates

General linear model

SPMs

**Multiple Comparisons**

# Statistical Parametric Mapping

➢ **From voxel-wise models to functional models**

➢ **Multiscale-adaptive estimation and inference procedures**

➢ **Wild-bootstrap methods to correct for multiple comparisons**

**Image=$f$(B(age, gene, race, disease, others), device, acquisition, noises)**

## Parametric, Semiparametric and Nonparametric Models for OOD analyses



$$= g(x, \theta, f) \oplus \varepsilon$$

$$x \in R^k, \theta \in \Theta \subset R^p, f \in F$$

$$g : R^k \times R^p \times F \to M$$

# Intrinsic Regression Models

- **Feature Methods:** Use some feature extraction functions to project random objects to Euclidean-valued variables.

- **Extrinsic Methods:** Ignore the fact that manifold-valued data are in a nonlinear space and then directly apply classical multivariate regression.

- **Intrinsic Methods:** few parametric models for manifold-valued data .

**Extrinsic Methods**

**Intrinsic Methods**



Dryden, I.L., Koloydenko, A. and Zhou, D. (2008).

# Intrinsic Regression Models

## Geodesic Link Function



$$f(x) = \text{Exp}(p, xv)$$

$$f(x_i) = \text{Exp}(p_1, (x_i - \bar{x})v_1)$$

**Single-center**

**Fletcher (2013)**
**Maxwell et al. (2014)**

$$g(x_i, q) = q_0 + x_i q_1 = q_0 + \bar{x} q_1 + (x_i - \bar{x}) q_1$$

$$f(x_i) = \text{Exp}(p_1, (x_i - \bar{x})v_1)$$

$$f(x_i) = \text{Exp}(p_1, \overset{\circ}{\underset{k}{a}} (x_{ik} - \bar{x}_k)v_k)$$

# Intrinsic Regression Models

**Residual**

$M$ **How to define residual?**

$T_D M$

$D$

$e_D$

$Y = Exp_D(e_D)$

**Inner product** $T_D M$

$$<< e_D, \tilde{e}_D >>$$

**Geodesic** $g_D(t; e_D)$

**Riemannian exponential maps**

$$Y = Exp_D(e_D) = g_D(1; e_D)$$

**Riemannian logarithm maps**

$$e_D = Log_D(Y) \, \widehat{Ì} \, B(0, r) \, Ì \, T_D M$$

*radius of injectivity*

# Intrinsic Regression Models

## Conditional Mean

### Riemannian logarithm maps

$$e(x,q,b) = \text{Log}_{m(x,q,b)}(Y) \hat{1} \ T_{m(x,q,b)}M$$

### Conditional Moment Model

$$E[e(x,q,b)\,|\,x] = E[\text{Log}_{m(x,q,b)}(Y)\,|\,x] = 0$$

$$T_{m(x,q,b)}M$$

$$M$$

$$m(x,q,b)$$

$$e(x,q,b)$$

$$Y = Exp_{m(x,q,b)}(e(x,b))$$

$$g_{m(x,q,b)}(t;e(x,q,b))$$

$$x$$

**Cornea, E**., **_Zhu, H.T.,_** Kim, P. and Ibrahim, J. G. Intrinsic regression model for data in Riemannian symmetric space. *JRSS, Series B*, 79, 463-482, 2017.

# Imaging Genetics of Brain Disorders

Most major brain disorders (like AD) are **heritable complex traits/diseases**

Together 50%-70% of AD risk
75%-90% of ADHD risk
60%-85% of Schizophrenia risk
~80% of Autism Spectrum Disorder (ASD) risk

Complex traits/diseases (many genes, environmental factors, complex functional mechanism)

Genetic signals are non-spare and weak:
Need large sample size to detect weak signals

Many genes contribute to the risk of AD
(polygenic genetic architecture)
(small but nonzero contribution)

# IG: Reproducibility and Heritability

# Brain- Heart Imaging Genetics Knowledge Portal



**Brain Imaging Genetics Knowledge Portal (BIG-KP)**



**Heart Imaging Genetics Knowledge Portal (Heart-KP)**

Aim to build the best knowledge database of neuroimaging genetics

# It's just a beginning

**Publications (2018+)**

Heart-brain connections: Phenotypic and genetic insights from magnetic resonance images. *Science* 380, abn6598 (2023). LINK.

Genetic influences on the shape of brain ventricular and subcortical structures (2022). medRxiv, 22270601 LINK

Common variants contribute to intrinsic human brain function networks (2022). *Nature Genetics*. LINK.

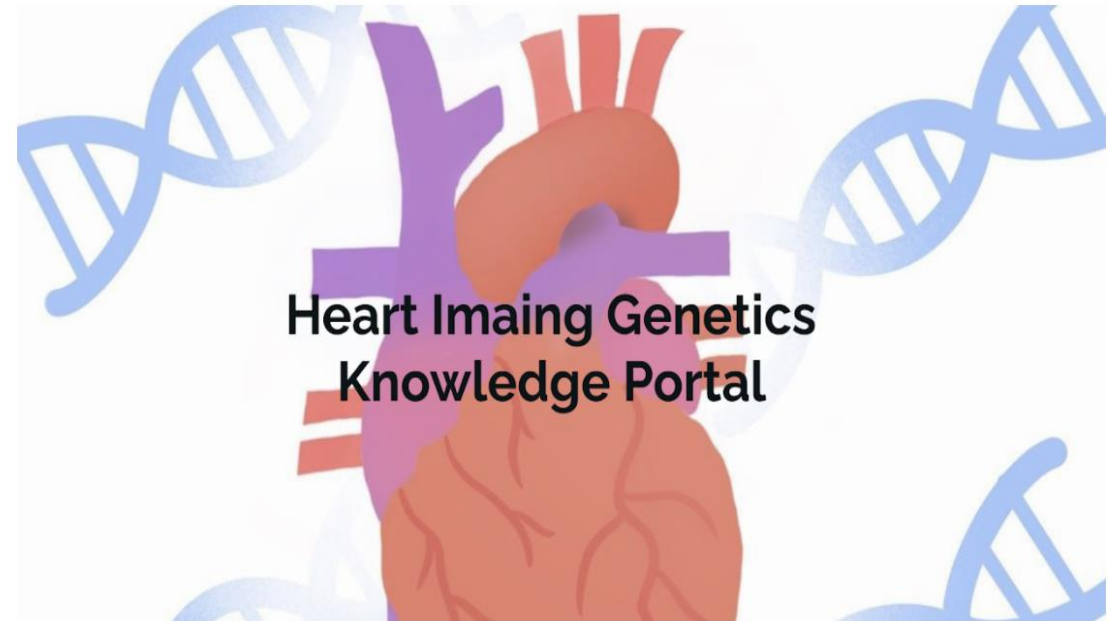Genetic influences on the intrinsic and extrinsic functional organizations of the cerebral cortex (2021). **medRxiv**, 21261187. LINK

Common genetic variation influencing human white matter microstructure (2021). *Science*, 372-6548. LINK

Transcriptome-wide association analysis of brain structures yields insights into pleiotropy with complex neuropsychiatric traits (2021). *Nature Communications*, 842872. LINK

Genome-wide association analysis of 19,629 individuals identifies variants influencing regional brain volumes and refines their genetic co-architecture with cognitive and mental health traits (2019). *Nature Genetics*, 51(11), 1637-1644. LINK

Large-scale GWAS reveals genetic architecture of brain white matter microstructure and genetic overlap with ... traits (n= 17,706) (2019). *Molecular Psychiatry*, in press. LINK

Heritability of regional brain volumes in large-scale neuroimaging and genetic studies (2018). *Cerebral Cortex*, 29(7), 2904-2914. LINK

Hundreds of associated genetic variants for 2100+ neuroimaging traits across three modalities: (grey matter volume, white matter microstructure, resting-state functional connectivity+rfMRI, task fMRI, shape, heart )

We make our research results publicly available by building the following resources.

If you are interested in other summary-level data from our analyses or have any questions or comments, feel free to contact **Bingxin Zhao (bingxin@purdue.edu)** or **Hongtu Zhu (htzhu@email.unc.edu)**.
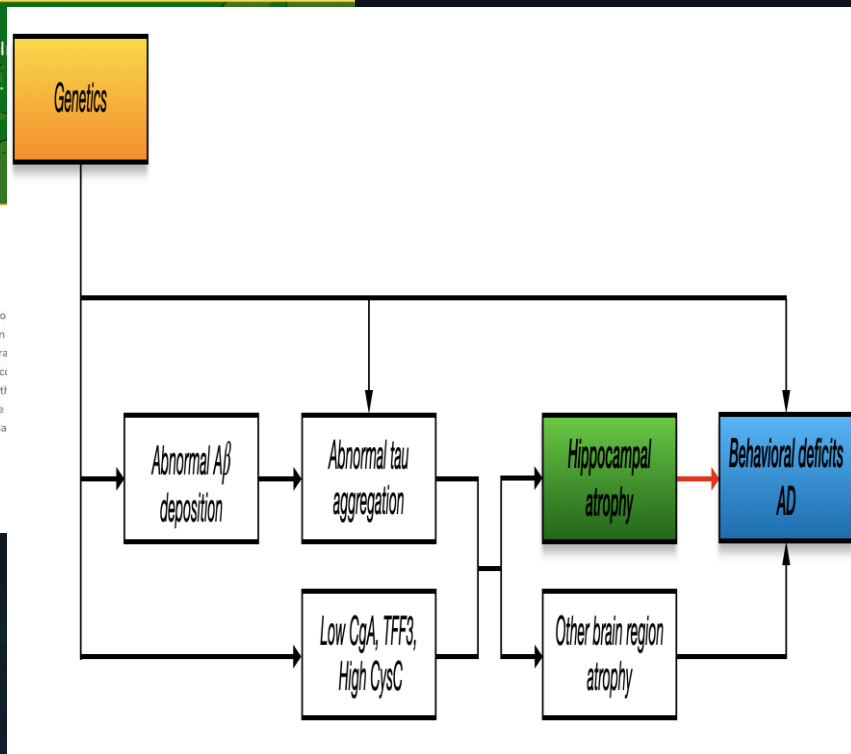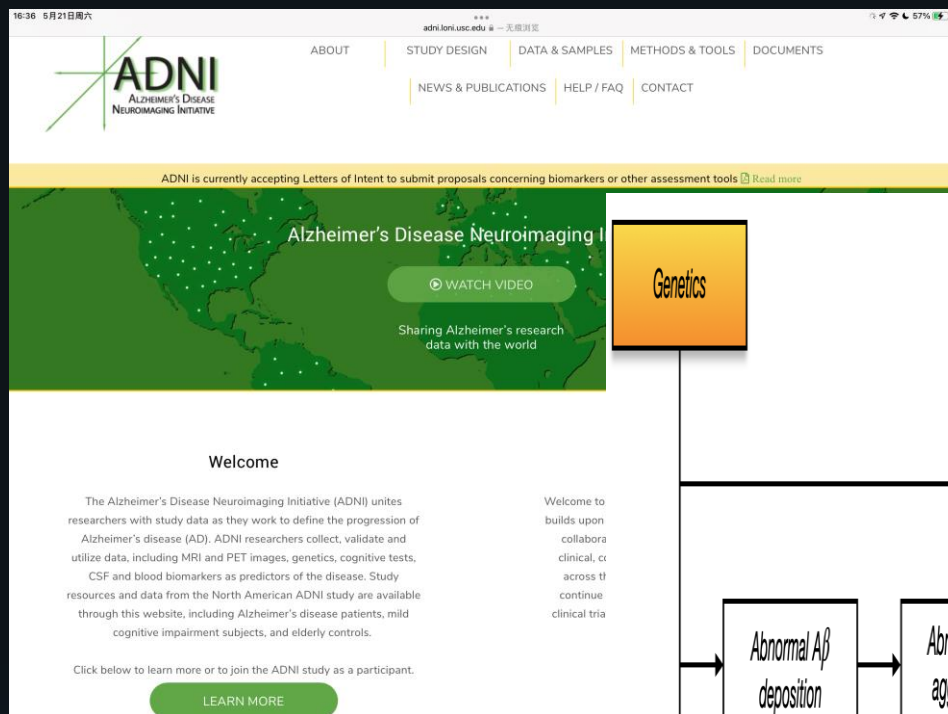
## 1. Imaging Genetics Online Server

We build a GWAS browser using the PheWeb tool to explore GWAS results for massive functional, structural, and diffusion neuroimaging traits. Currently, we support GWAS results of 2104 traits trained in the UKB British cohort (n~34,000), including

1. 635 ENIGMA-DTI parameters of brain white matter (diffusion MRI)
2. 376 ANTS regional brain volumes (structural MRI)
3. 191 ICA-based functional MRI traits (rs-fMRI(ICA))
4. 309 parcellation-based functional MRI (task/rs-fMRI(Glasser360))

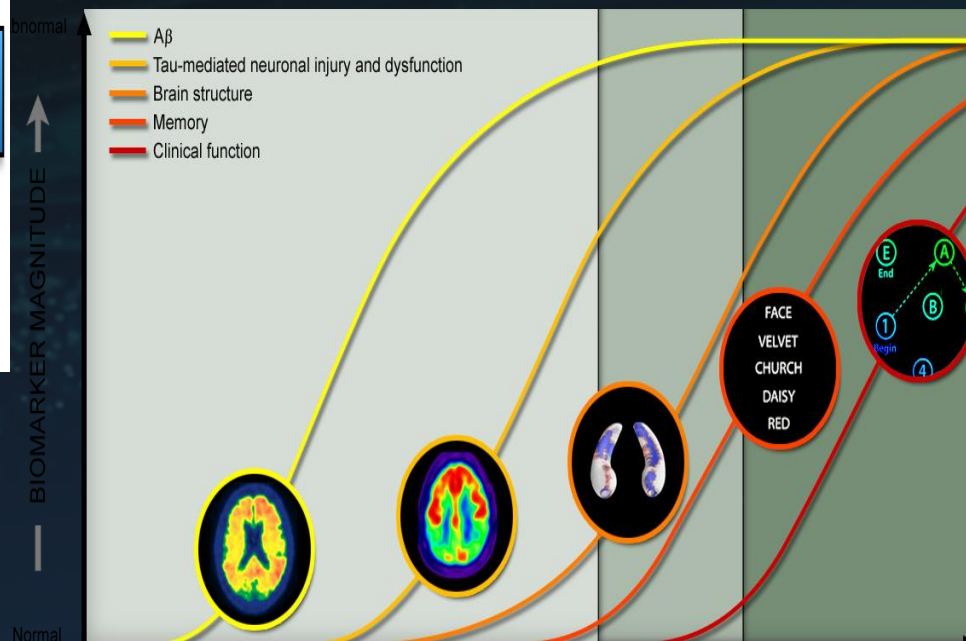# Genetics discovery in human brain by big data integration

# Alzheimer's Disease Neuroimaging Initiative



**For a heterogeneous, clinically defined disorder, the endophenotype is 'closer to the underlying biology,'**
- **Increasing the power of genetic search**
- **Being informative about disorder risk.**
- **Providing mechanistic connections linking genetic variation to behavioral measures.**

2004-now

# Model Setup

**Outcome generating model**

$$Y_i = \sum_{l=1}^{s} x_{il} \beta_l + < \boldsymbol{Z_i}, \boldsymbol{B} > + \epsilon_i$$

**Exposure generating model**

$$\boldsymbol{Z_i} = \sum_{l=1}^{s} x_{il} * \boldsymbol{C_l} + \boldsymbol{E_i}$$

$\boldsymbol{B}$ is the main parameter of interest, representing the association between the 2D imaging exposure $Z_i$ and the behavioral outcome $Y_i$, $\beta_l$ represents the association between the l−th observed covariate $x_{il}$ and the behavioral outcome $Y_i$, and $\epsilon_i$ and $E_i$ are random errors that may be correlated. The symbol "$*$" denotes element-wise multiplication.

Ye, Wang, Kong, and Zhu (2022). Mapping the Genetic-Imaging-Clinical Pathway with Applications to Alzheimer's Disease. JASA, in press.

# Marginal Screening

**Fit:**

$$Y_i = x_{il}\beta_l + \epsilon_i$$

**Obtain:**

$$\hat{\beta}_l^M = n^{-1} \sum_{i=1}^n x_{il} Y_i$$

**Problem!!! (plugging exposure model into outcome model)**

**Outcome generating model** $Y_i = \sum_{l=1}^s x_{il} \beta_l + < Z_i, B > + \epsilon_i$

**Exposure generating model** $Z_i = \sum_{l=1}^s x_{il} * C_l + E_i$

Obtain:

$$Y_i = \sum_{l=1}^s x_{il} (\beta_l + < C_l, B >) + < E_i, B > + \epsilon_i$$

Miss a portion of confounders when $\beta_l$ and $< C_l, B >$ are of similar magnitude but opposite sign.

# Joint Screening (proposed)

Marginal screening:

$$\mathbf{Z}_i = \sum_{l=1}^{s} x_{il} * \mathbf{C}_l + \mathrm{E}_i$$

Obtain (*Kong, An, Zhang and Zhu, 2020*):

$$\widehat{\boldsymbol{C}}_l^M = n^{-1} \sum_{i=1}^{n} x_{il} * \mathbf{Z}_i \in \mathbb{R}^{p \times q}$$

$$\widehat{\mathcal{M}}_1^* = \left\{ 1 \leq I \leq s : \left| \widehat{\beta_l^M} \right| \geq \gamma_{1,n} \right\}$$

$$\widehat{\mathcal{M}}_2 = \left\{ 1 \leq I \leq s : \| \widehat{\boldsymbol{C}}_l^M \|_{op} \geq \gamma_{2,n} \right\}$$

$$\mathcal{C} = \{ l \in \mathcal{A} \mid \beta_l \neq 0 \text{ and } \boldsymbol{C}_l \neq 0 \},$$
$$\mathcal{P} = \{ l \in \mathcal{A} \mid \beta_l \neq 0 \text{ and } \boldsymbol{C}_l = 0 \},$$
$$\mathcal{I} = \{ l \in \mathcal{A} \mid \beta_l = 0 \text{ and } \boldsymbol{C}_l \neq 0 \},$$
$$\mathcal{S} = \{ l \in \mathcal{A} \mid \beta_l = 0 \text{ and } \boldsymbol{C}_l = 0 \}.$$

**Select submodel:** $\widehat{\mathcal{M}} = \widehat{\mathcal{M}}_1^* \cup \widehat{\mathcal{M}}_2.$ (Union)

**Alternative choices (both worse):** $\widehat{\mathcal{M}}_1^*$ (outcome) or $\widehat{\mathcal{M}}_1^* \cap \widehat{\mathcal{M}}_2$ (Outcome).

# Estimation (proposed)

**Minimize**:

$$\frac{1}{2}\sum_{i=1}^{n}\left(Y_i - \langle \mathbf{Z}_i, \mathrm{B} \rangle - \sum_{l\in\widehat{\mathcal{M}}} X_{il}\beta_l\right)^2 + \lambda_{1,n}\sum_{l\in\widehat{\mathcal{M}}}|\beta_l| + \lambda_{2,n}\parallel \mathrm{B} \parallel_*$$

where $\parallel \mathrm{B} \parallel_* = \sum_k \sigma_k(\mathrm{B})$ .

L1 penalty, exclude instrumental and irrelevant variables.

Nuclear penalty, low-rank estimation of B.

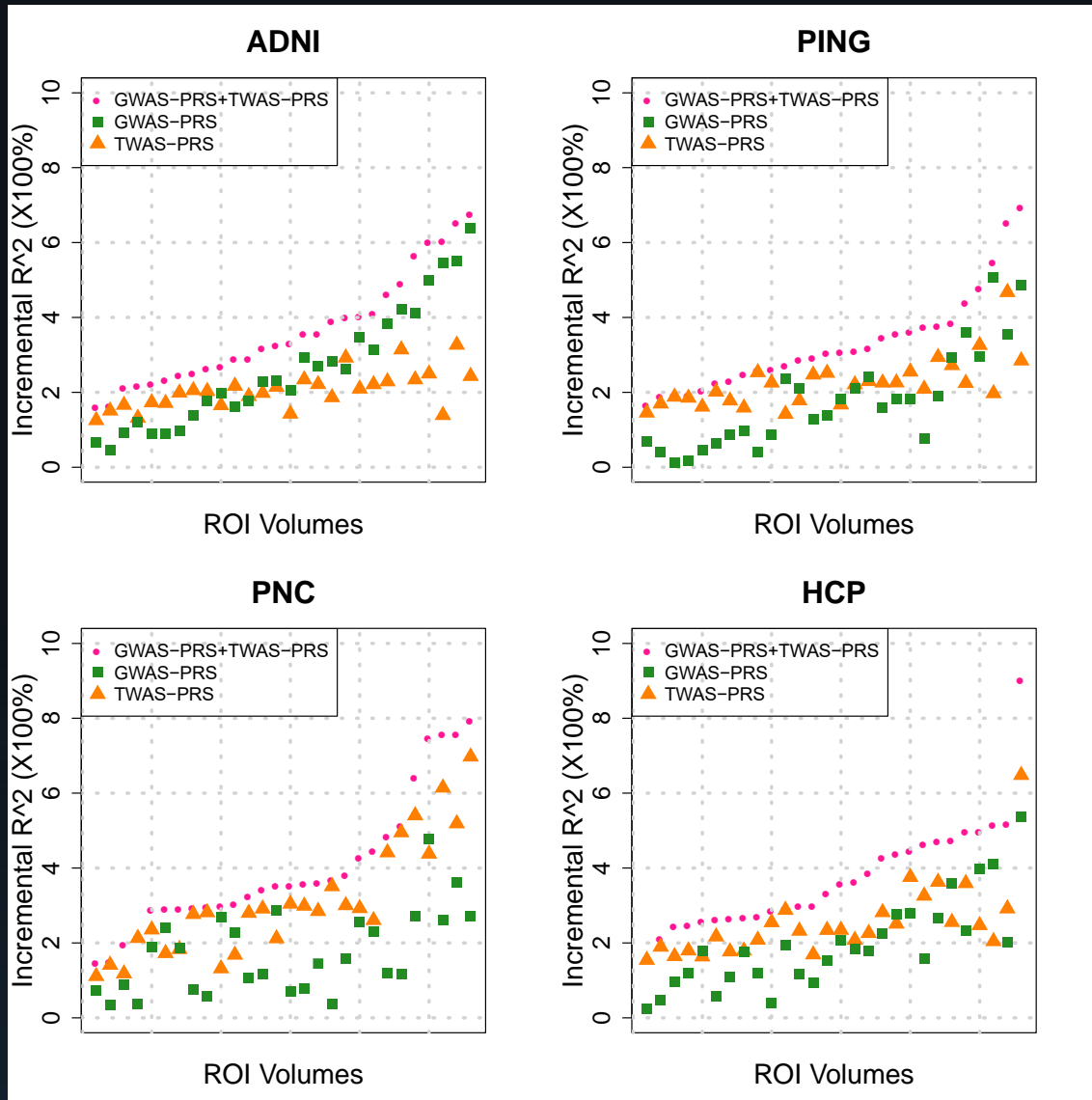Estimated effect size of imaging exposure $z$,

$$\hat{\mu}(z) = \langle z, \hat{B} \rangle$$

$$\mathcal{C} = \{l \in \mathcal{A} \mid \beta_l \neq 0 \text{ and } \boldsymbol{C}_l \neq 0\},$$
$$\mathcal{P} = \{l \in \mathcal{A} \mid \beta_l \neq 0 \text{ and } \boldsymbol{C}_l = 0\},$$
$$\mathcal{I} = \{l \in \mathcal{A} \mid \beta_l = 0 \text{ and } \boldsymbol{C}_l \neq 0\},$$
$$\mathcal{S} = \{l \in \mathcal{A} \mid \beta_l = 0 \text{ and } \boldsymbol{C}_l = 0\}.$$
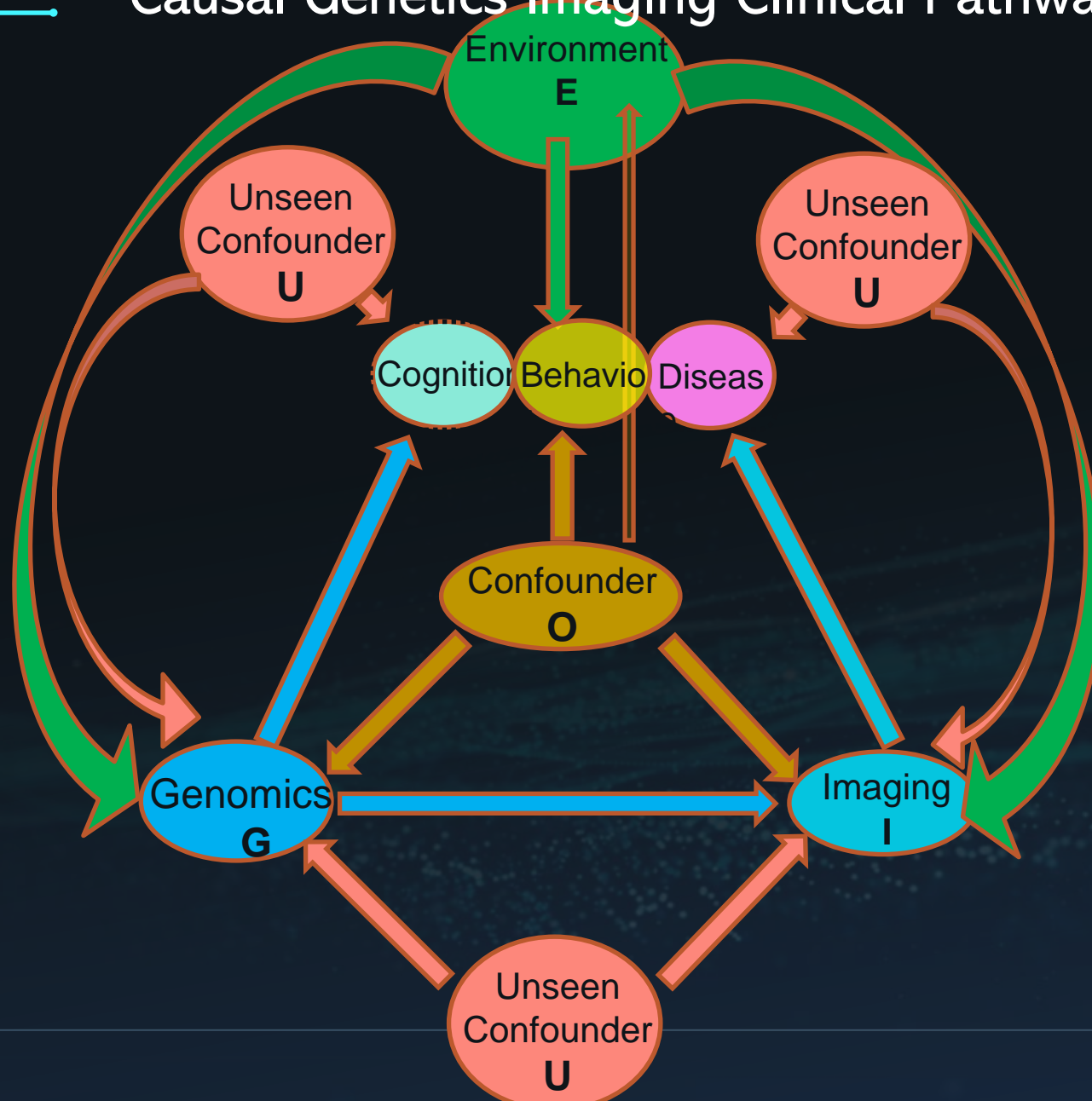
# Predictive Analysis



Gene expression-informed gene-level PRS + GWAS PRS has higher prediction accuracy

Construct gene-level PRS (polygenic risk scores) by leveraging gene expression reference panels (e.g., GTEx) in TWAS
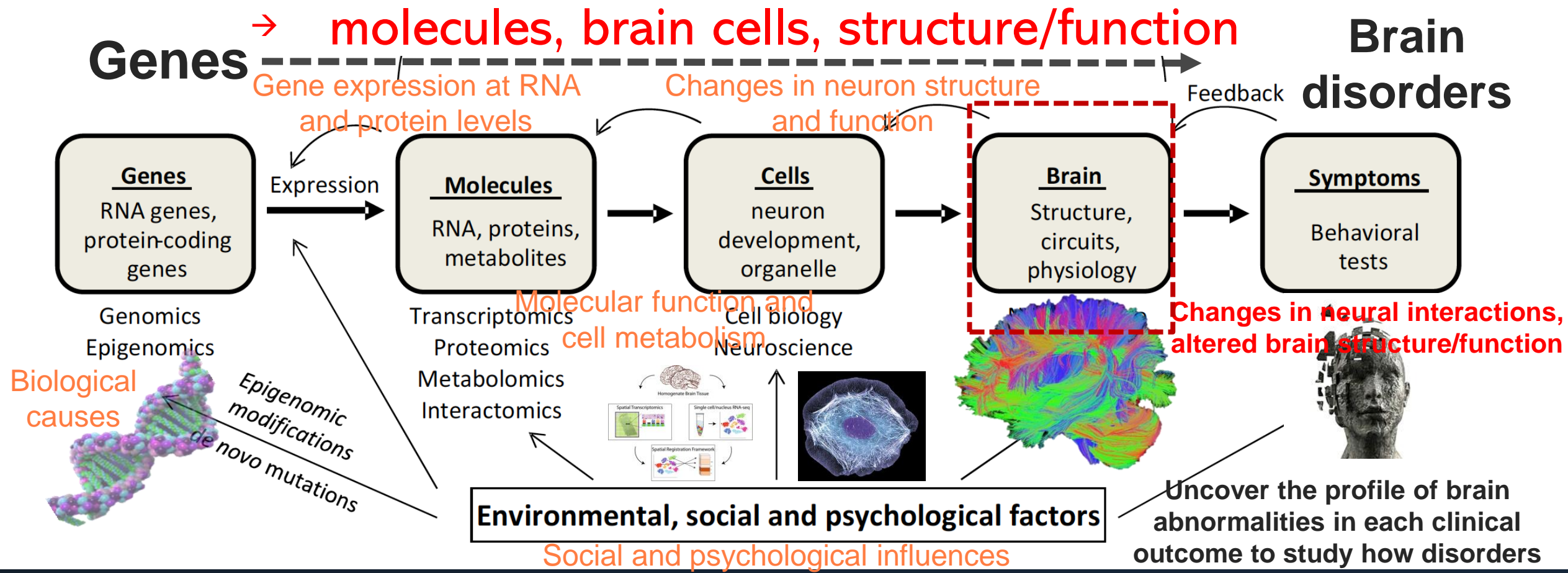
Causal Genetics Imaging Clinical Pathway

# Brain Imaging Genetics Paradigm

Neuroimaging: an important component to help understand the complex biological pathways of brain disorders



**Genes** → **molecules, brain cells, structure/function** → **Brain disorders**

Gene expression at RNA and protein levels

Changes in neuron structure and function

Feedback

| Genes | Expression | Molecules | Cells | Brain | Symptoms |
|---|---|---|---|---|---|
| RNA genes, protein-coding genes | | RNA, proteins, metabolites | neuron development, organelle | Structure, circuits, physiology | Behavioral tests |

Genomics
Epigenomics

Biological causes

Epigenomic modifications

de novo mutations

Molecular function and cell metabolism

Transcriptomics
Proteomics
Metabolomics
Interactomics

Cell biology
Neuroscience

**Changes in neural interactions, altered brain structure/function**

**Environmental, social and psychological factors**

Social and psychological influences

**Uncover the profile of brain abnormalities in each clinical outcome to study how disorders develop**

# Challenges

❖ **The complexity of those large-scale neuroimaging-related data sets is too high for most research teams in both academia and industry.**

❖ **It is very difficulty to appropriately process data across different domains with high quality, while controlling for potential bias introduced during the preprocessing stage.**

❖ **It remains uncertain as to how to appropriately integrate data across different domains obtained from different studies and cohorts with possible different study designs for unbiased data integration.**

❖ **It remains unclear how to appropriately and efficiently analyze neuroimaging related data sets with multiple Vs (e.g., Volume, Velocity, Variety and Veracity), while ensuring algorithmic fairness.**

# Statistical Learning Methods for NDA



**(Zhu, Li & Zhao, 2023)**

# Acknowledgement

**Brain Imaging Genetics Knowledge Portal (BIG-KP)**

Genetics Discoveries in Human Brain by Big Data Integration

bigkp.org