

Linear Mixed Models minimise false positive rate and enhance precision of mass-univariate vertex-wise analyses of grey-matter

Baptiste Couvy-Duchesne,
Futao Zhang, Kathryn E. Kemper,
Julia Sidorenko, Naomi R. Wray,
Peter M. Visscher, Olivier Colliot,
Jian Yang

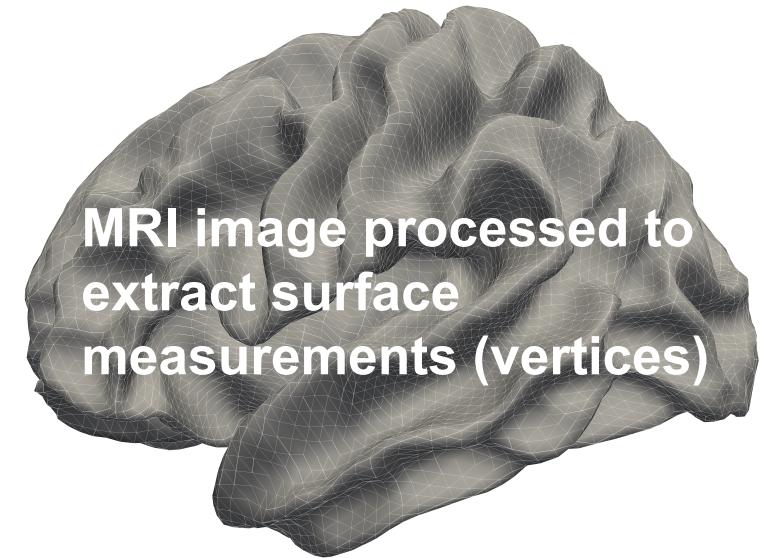
b.couvyduchesne@uq.edu.au /
baptiste.couvy@icm-institute.org



Mass-univariate vertex-wise analyses

Disease / trait

(e.g. Alzheimer's,
cognition score,
smoking status...)



Objective: find brain regions associated with phenotype

Context: “big-data” problem with more vertices/feature than participants

Current approach: test association with each vertex in turn (using **Generalised Linear Models: GLM**)

Rationale

- Vertices are correlated with each other due to factors unrelated to the phenotype of interest (e.g. geometry, MRI artifacts, age, sex)
 - Vertices correlated with truly associated ones may also reach significance
 - **PROBLEM: confounded association, pointing to “wrong” brain regions and redundant**
- **Linear Mixed Models (LMM)** can control for factors that cause correlation between features, leading to more robust association results (evidence from the OMICs literature)



NIH Public Access Author Manuscript

Nat Genet. Author manuscript; available in PMC 2014 August 01.

Published in final edited form as:

Nat Genet. 2014 February ; 46(2): 100–106. doi:10.1038/ng.2876.

Advantages and pitfalls in the application of mixed model association methods

Jian Yang^{1,2,*}, Noah A. Zaitlen^{3,*}, Michael E. Goddard^{4,**}, Peter M. Visscher^{1,2,**}, and Alkes L. Price^{5,6,7,**}

¹University of Queensland Diamantina Institute, University of Queensland, Princess Alexandra Hospital, Brisbane, Queensland, Australia

Jian Yang , Weinan Chen , Zihong Zhu , Qian Zhang¹, Marta F. Nabais^{1,2}, Ting Qi¹, Ian J. Deary³, Naomi R. Wray^{1,4}, Peter M. Visscher^{1,4}, Allan F. McRae¹ and Jian Yang^{1,4,5*} 

Genome Biology

Open Access

ic-data-based complex



Aims

Compare state of the art (GLM) to Linear Mixed Models (LMM) for association studies of the brain

- False positive rate
- Statistical power
- Precision (size of associated regions/clusters)
- Prediction accuracy

- 1) Using phenotypes simulated from real MRI images
- 2) Using real phenotypes (age, sex, BMI, fluid IQ and smoking status)



T1w MRI image

DATA:

Processed T1w and T2 FLAIR MRI images

Vertices measure grey-matter structure (thickness and surface area)

UK Biobank: 10,102 volunteers from the UK, aged 44-80

For prediction and validation:

OASIS3: 732 adults some with Alzheimer's disease, aged 42-95

UK Biobank: another 4,942 participants from the UK

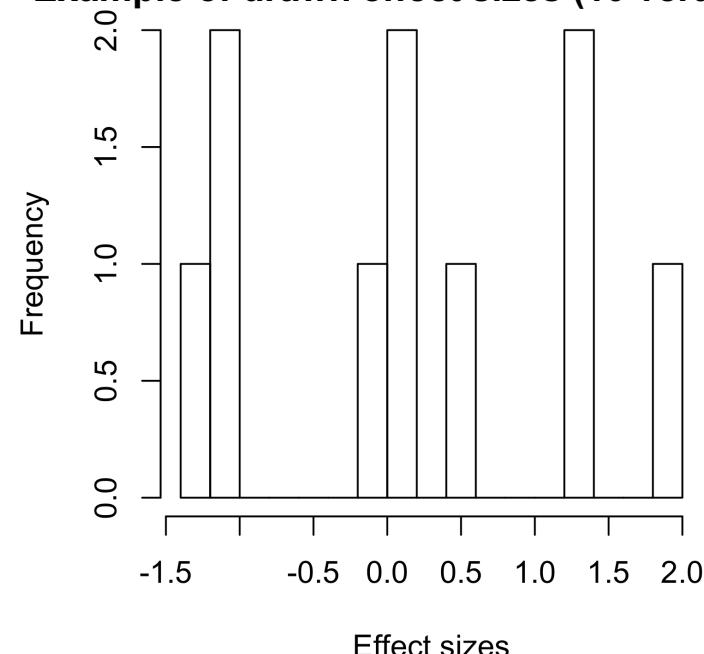


Simulation scenarios

Phenotype simulation

- Randomly select associated vertices
 - draw effect sizes (b_i) from a normal distribution
 - Set the total association between vertices and phenotype (R^2 a.k.a. morphometricity)
 - Construct simulated phenotype y as
- $y = \sum(x_i b_i) + \epsilon$ with:
- x_i vector of standardised vertex measurement, at the i -th associated vertex
- $\epsilon \sim N(0, \text{var}(\sum(x_i b_i))(1/R^2-1))$

Example of drawn effect sizes (10 vertices)



Scenario 1:

10 associated vertices
Total $R^2=0.2$

\Rightarrow Average vertex association
 $r^2=0.02$

Scenario 2:

100 associated vertices
Total $R^2=0.5$

$\Rightarrow r^2=0.005$

Scenario 3:

1000 associated vertices
Total $R^2=0.4$

$\Rightarrow r^2=0.0004$

For each scenario, we simulated 100 phenotypes

Models considered

GLM no covariates:

$$\text{pheno} = u + b^*\text{vertex} + \varepsilon$$

GLM with age, sex and ICV:

$$\text{pheno} = u + b_1^*\text{age} + b_2^*\text{sex} + b_3^*\text{ICV} + b^*\text{vertex} + \varepsilon$$

GLM with 5 principal components:

$$\text{pheno} = u + b_1^*\text{PC1} + \dots + b_5^*\text{PC5} + b^*\text{vertex} + \varepsilon$$

GLM with 10 principal components:

$$\text{pheno} = u + b_1^*\text{PC1} + \dots + b_{10}^*\text{PC10} + b^*\text{vertex} + \varepsilon$$

GLM with 10 modality specific Principal components: $\text{pheno} = u + b_1^*\text{PC1}_j + \dots + b_{10}^*\text{PC10}_j + b^*\text{vertex}_j + \varepsilon$
 With j the modality (cortical thickness, area, subcortical thickness or area) of the vertex

LMM (global BRM)

$$\text{pheno} = u + b^*\text{vertex} + \beta + \varepsilon \\ \text{with } \beta \sim N(0, \mathbf{BRM} * \sigma^2)$$

LMM (modality specific BRM)

$$\text{pheno} = u + b^*\text{vertex}_j + \beta_1 + \beta_2 + \beta_3 + \beta_4 + \varepsilon \\ \text{with } \beta_j \sim N(0, \mathbf{BRM}_j * \sigma_j^2) \text{ the random effect corresponding to a single modality}$$

We apply **Bonferroni correction for multiple testing** in all the following

Simulation results



Inflation of ("null") test statistics

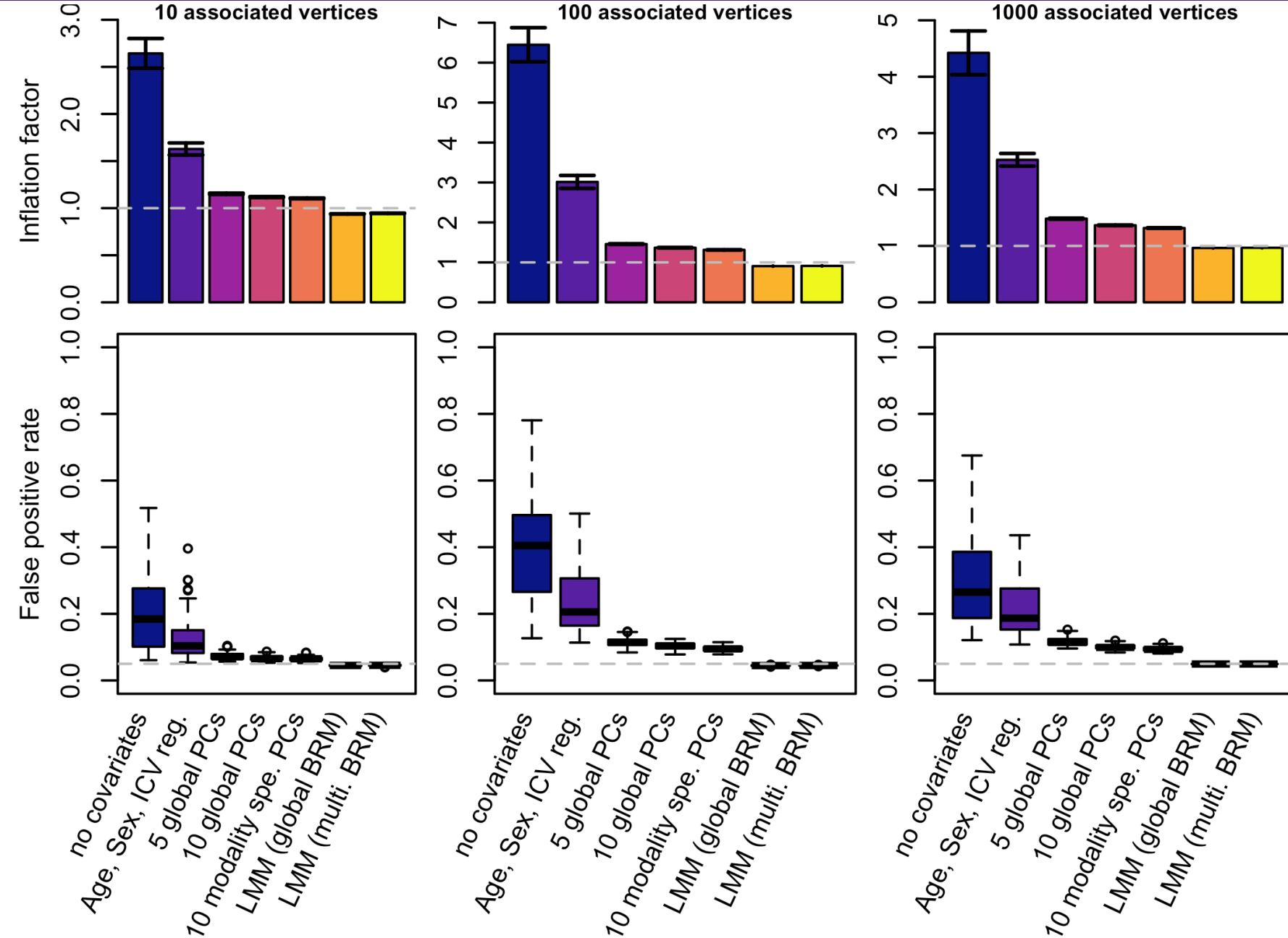
Compare observed and expected distribution of test statistics of non-associated vertices

Inflation factor: median of observed test statistic divided by the expected (under the null)

False positive rate: in each replicate the proportion of null vertices with $p < 0.05$

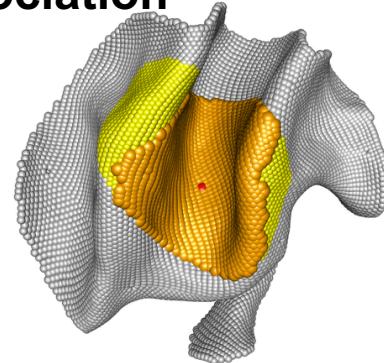
Inflation of statistics in GLM

- ⇒ Null vertices are correlated with true associations
- ⇒ Possible false positives



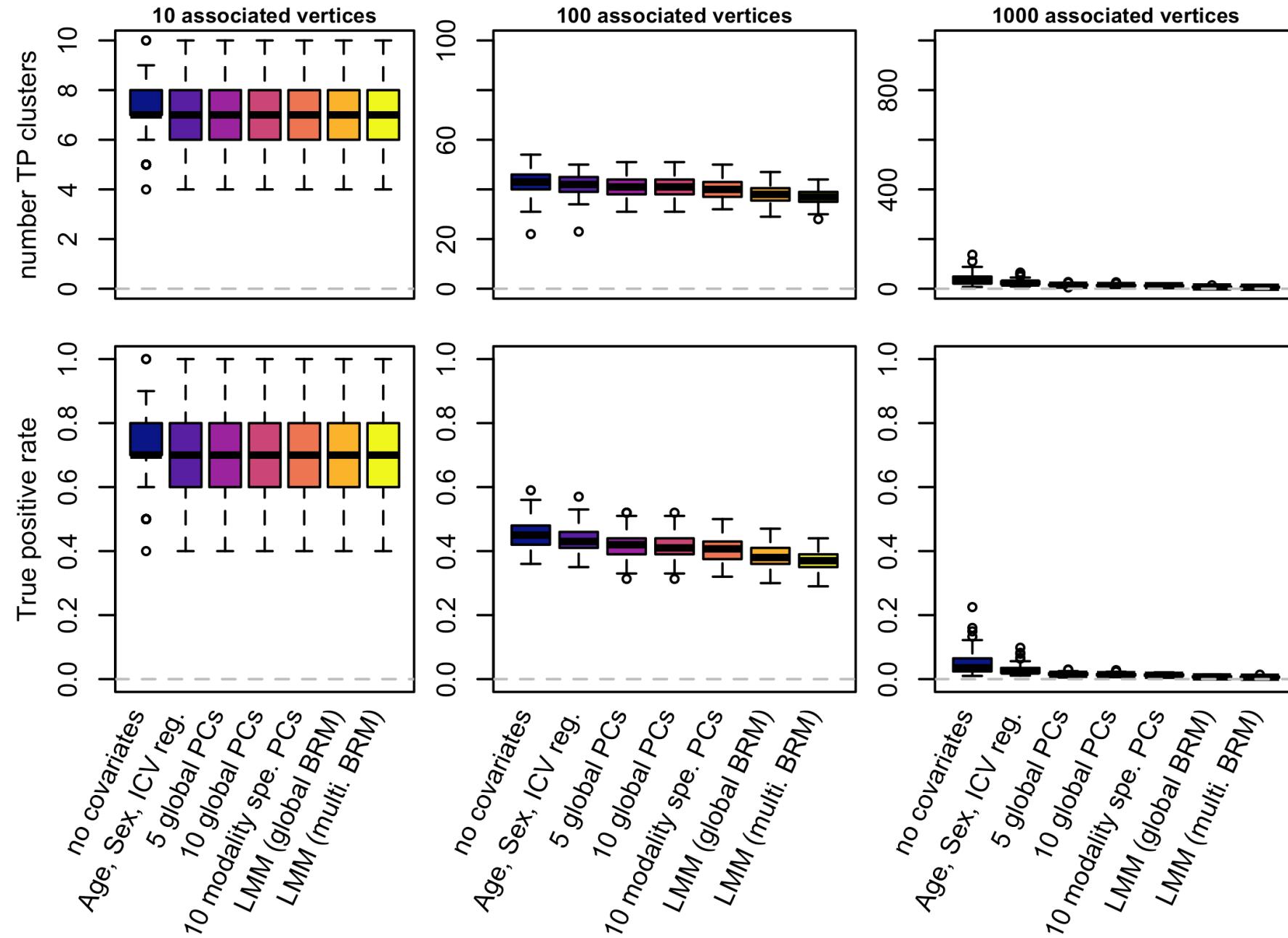
Statistical Power

Number of significant clusters that contain a true association



True positive rate: in each replicate, the proportion of truly associated vertices significant

Power slightly reduced for LMM compared to GLM



Precision

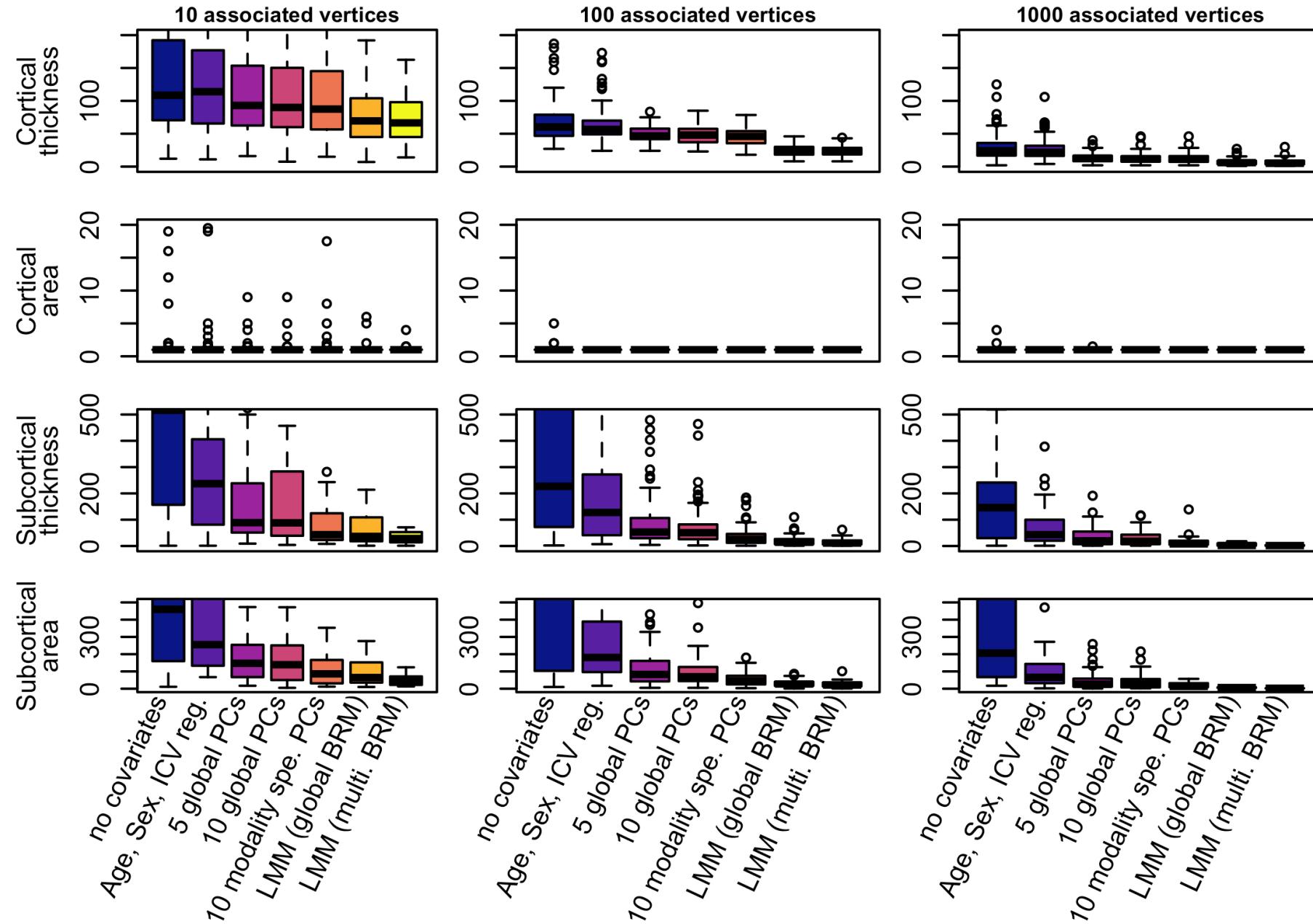
Median size of True Positive clusters
By modality

LMM yield true positive clusters much smaller than GLM

⇒ More precise identification of associated brain region

For cortical area more than half the TP clusters contain a single vertex

⇒ Suggests little local/long range correlations with those vertices



False positive rate

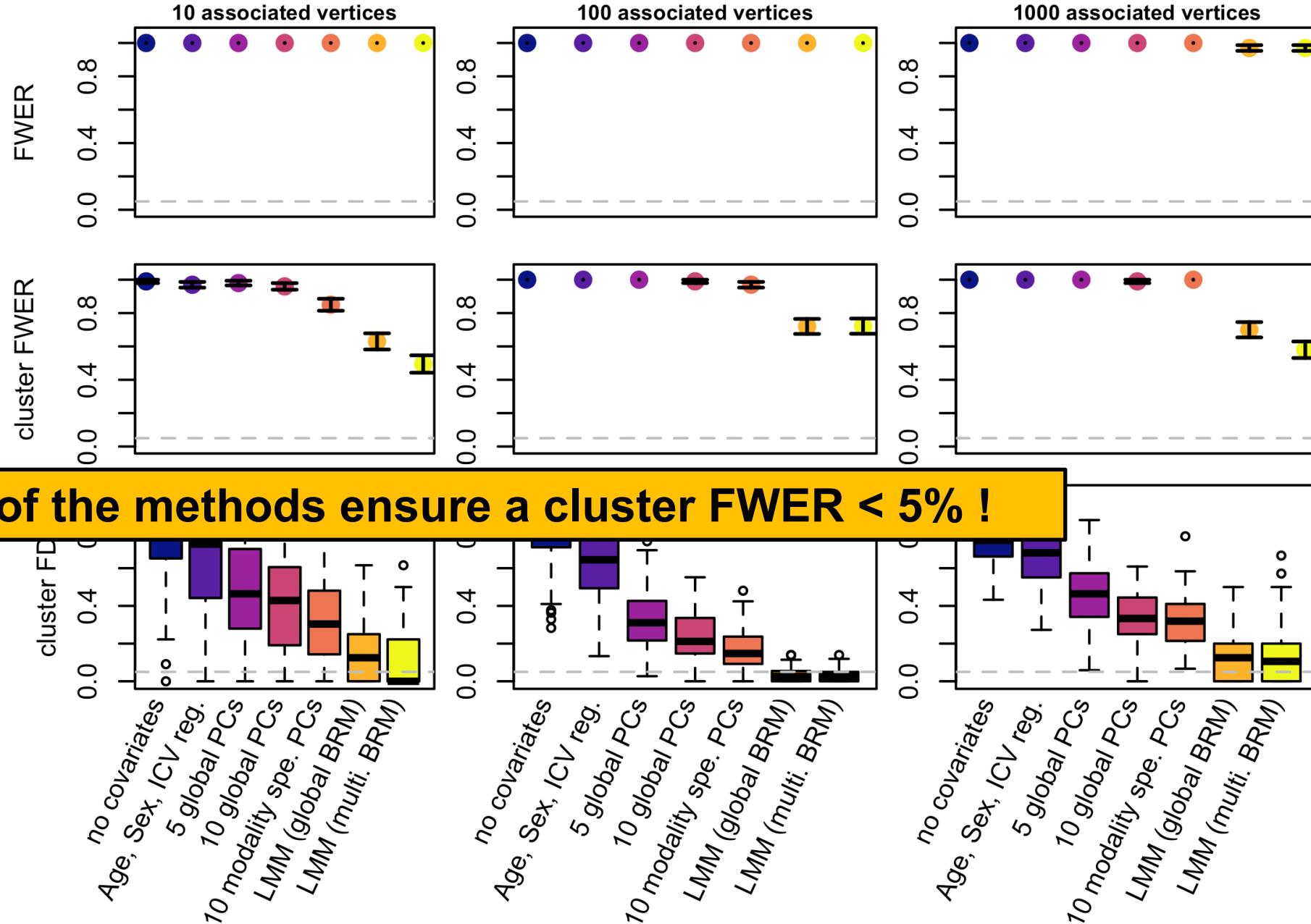
FWER (family wise error rate): proportion of replicates with at least 1 false positive vertex

Cluster FWER: proportion of replicates with at least 1 false positive cluster

Cluster FDR: proportion of FP clusters among significant clusters

All methods yield 1+ false positive vertex

LMM minimise probability of false positive clusters and proportion of false positive clusters



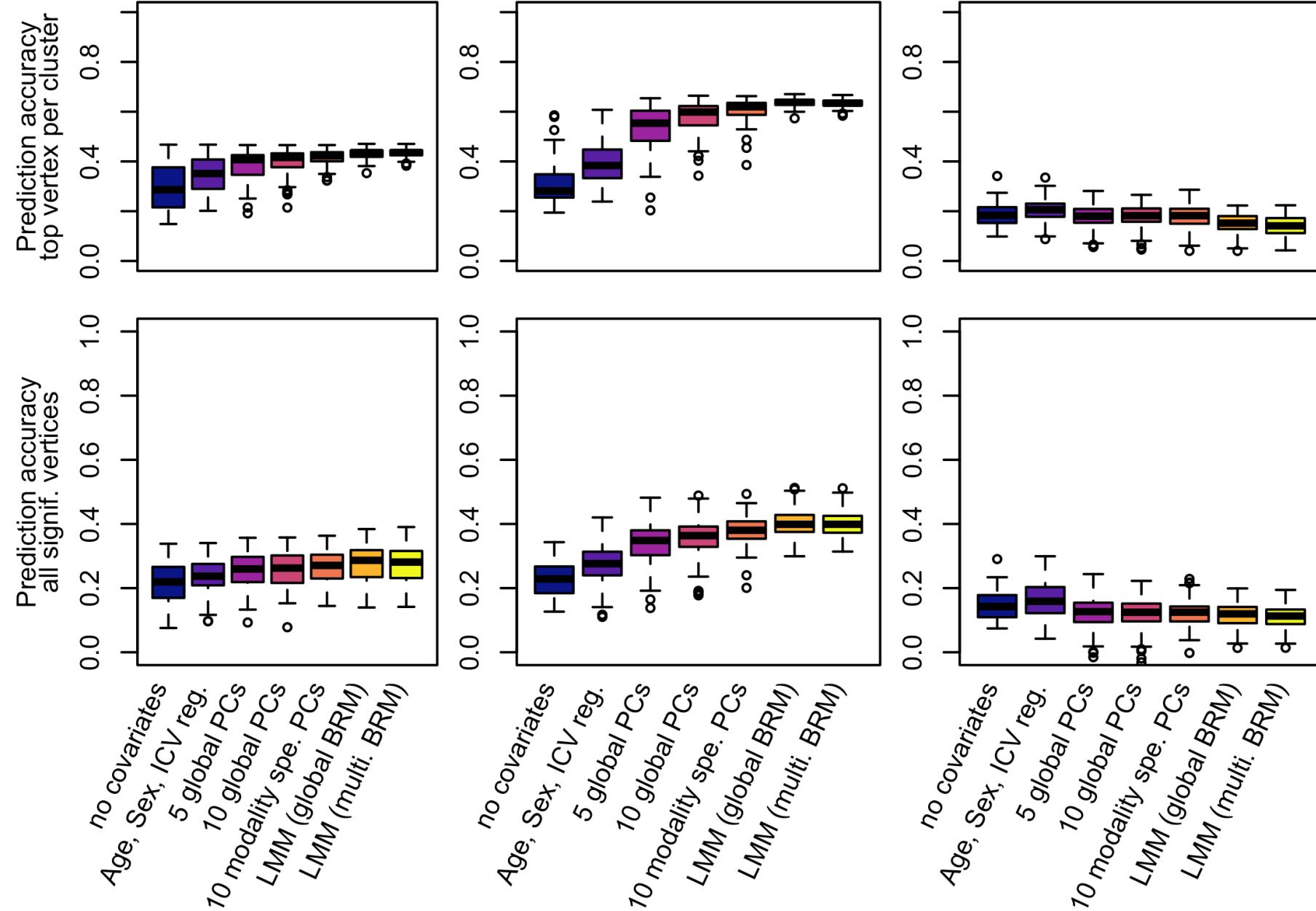
Prediction – into UKB left out sample

Prediction accuracy (correlation) from significant brain regions

- Using the top vertex (smallest pvalue) per significant cluster
- Using all significant vertices

LMM derived scores offer equivalent to superior prediction accuracy

GLM based prediction decreased by false positives and redundant associations



Summary of the simulations

**As in OMICs analyses,
LMM control for
unaccounted factors
that contribute to
correlation between
features which can
cause false positive**

**LMM offer a more
parsimonious and
conservative
characterisation of the
brain-trait associations**

	GLM	LMM
Inflation of “null” test statistics	Yes	No
Statistical power	Maximal	Slightly reduced
Precision	Poor	Good
FWER	Large (~1)	Large (~1)
Cluster FWER	Large (0.85 - 1)	Reduced (0.49 - 0.72)
Cluster FDR	High (0.17-0.78)	Small (0.03 - 0.16)
Computation	Quick (minutes)	Longer (hours)
Prediction accuracy	Better than chance	Equivalent to superior

Analyses of real phenotypes

Summary of significant associations

		BMI	Fluid Intelligence	Age	Sex	Smoking status
GLM no covariates	<i>N assoc. vertices</i>	10,977	24,809	136,348	355,127	712
	<i>N assoc. clusters</i>	232	640	970	714	34
	<i>Max cluster size</i>	862	2,030	22,358	130,651	116
GLM - Age, sex, ICV	<i>N assoc. vertices</i>	10,262	196	130,189	56,575	88
	<i>N assoc. clusters</i>	237	15	1,270	1,154	6
	<i>Max cluster size</i>	494	112	19,450	2,955	37
GLM - 10 global PCs	<i>N assoc. vertices</i>	8,565	30	16,772	26,912	91
	<i>N assoc. clusters</i>	174	5	297	492	6
	<i>Max cluster size</i>	518	9	894	1,782	43
Single random effect LMM	<i>N assoc. vertices</i>	11	0	47	27	0
	<i>N assoc. clusters</i>	5	0	8	6	0
	<i>Max cluster size</i>	5	NA	15	11	NA
	<i>Morphometricity (SE)</i>	0.59 (0.026)	0.16 (0.023)	0.91 (0.021)	0.99 (0.020)	0.15 (0.022)
Multiple random effect LMM	<i>N assoc. vertices</i>	0	0	0	9	0
	<i>N assoc. clusters</i>	0	0	0	1	0
	<i>Max cluster size</i>	NA	NA	NA	9	NA
	<i>Morphometricity (SE)</i>	0.51 (0.031)	0.17 (0.034)	0.83 (0.026)	1.06 (0.024)	0.12 (0.029)

As expected from simulations:

LMM yield fewer significant vertices or clusters than GLM, as well as smaller clusters

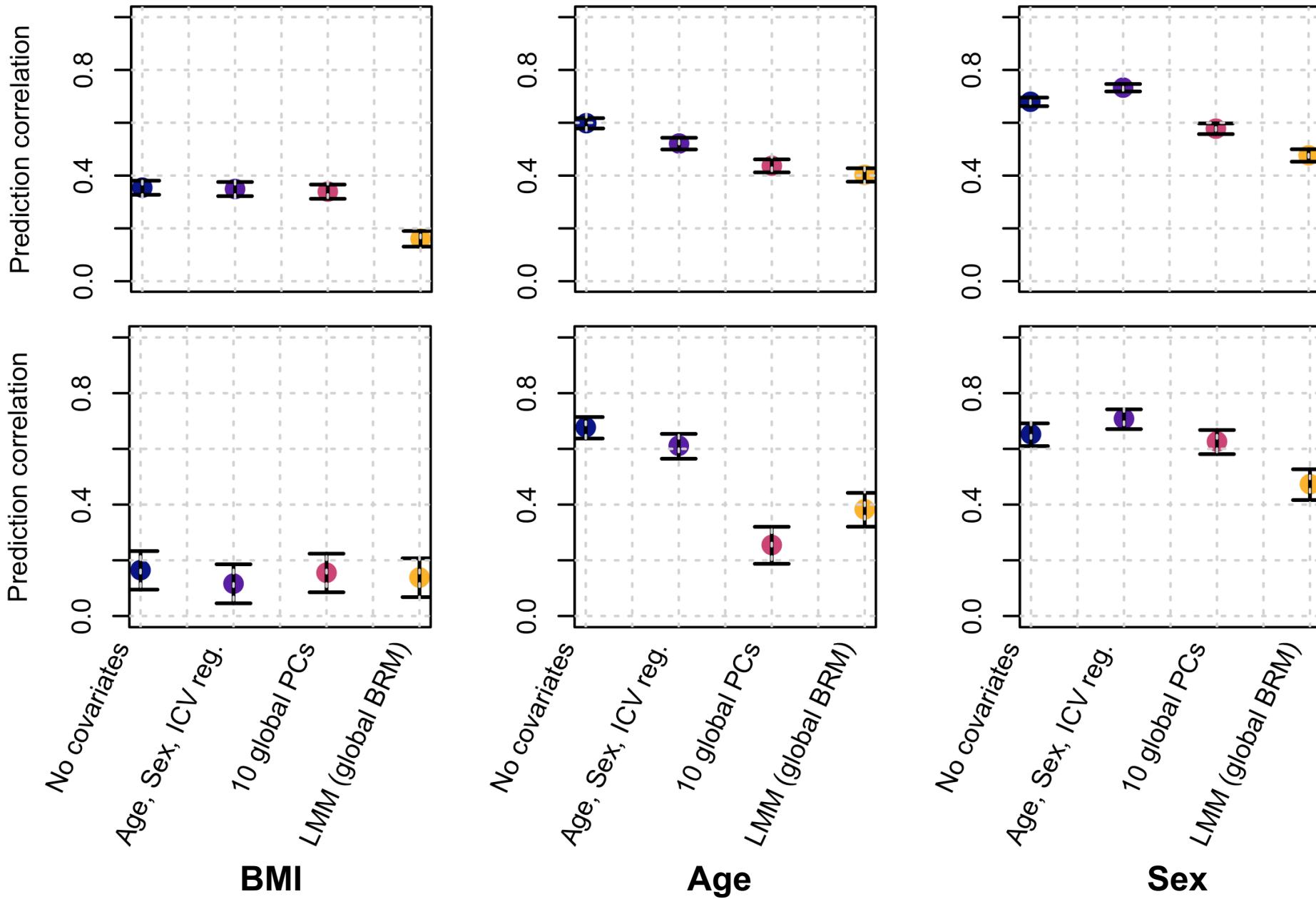
LMM with multiple random effects may suffer from limited power

Prediction from significant brain regions

Prediction into UKB left out sample

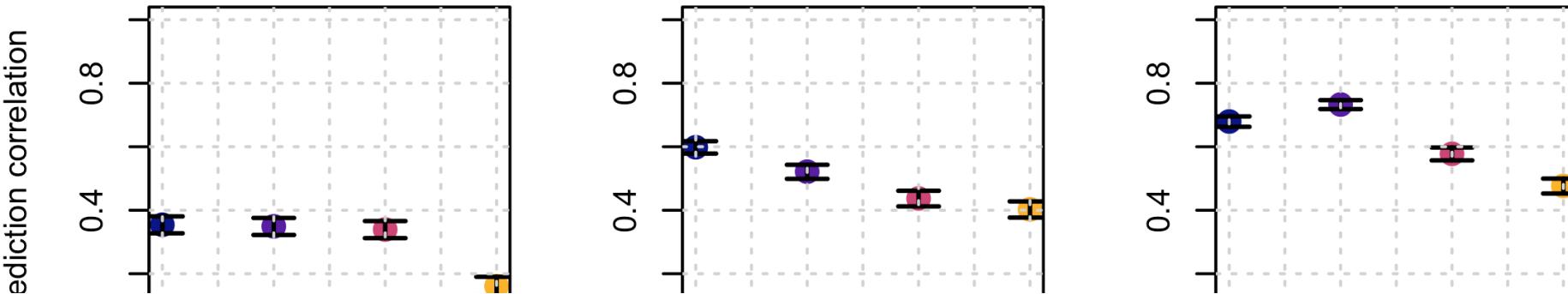
Prediction into OASIS3

BMI:
 GLM based prediction
 superior in UKB
But GLM fail to generalise in OASIS
 where the pattern of correlation between vertices may be different



Prediction from significant brain regions

Prediction into UKB left out sample

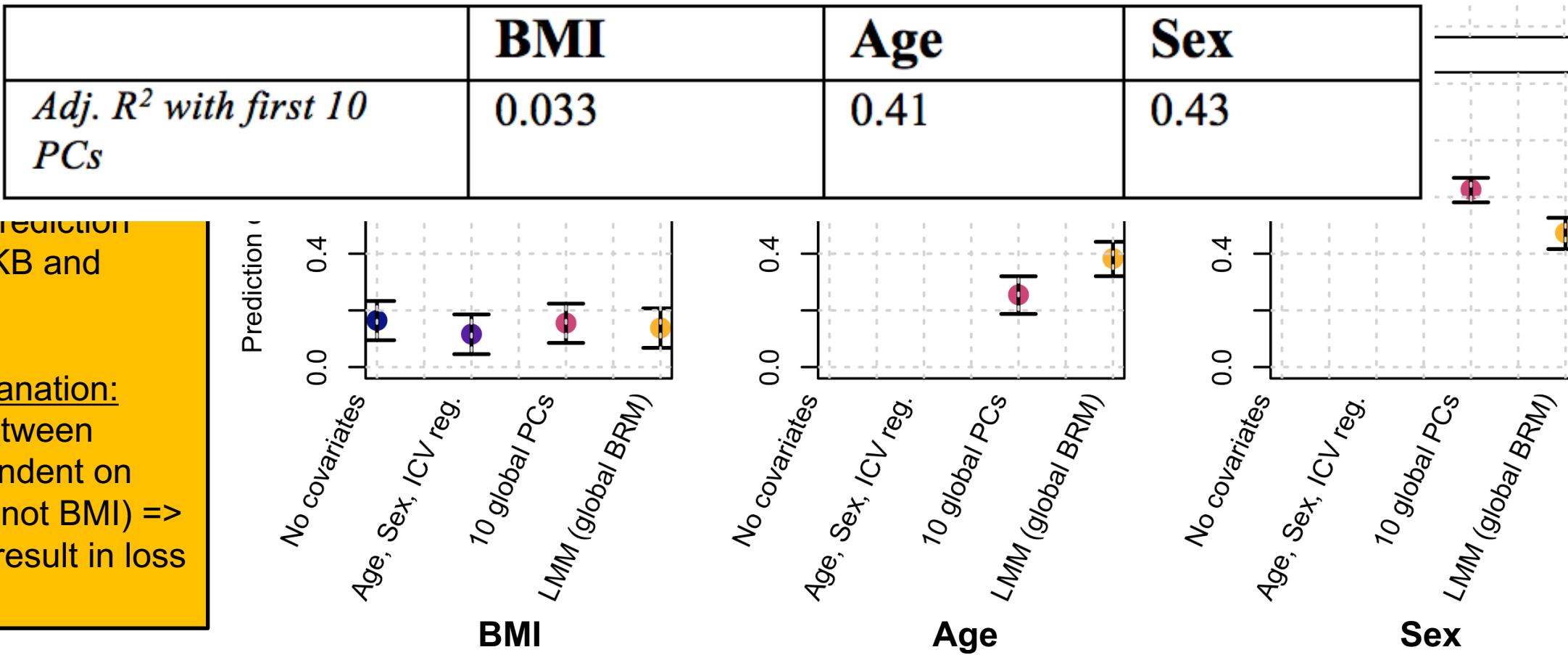


Prediction into OASIS

Age and sex

GLM based prediction superior in UKB and OASIS

Possible explanation:
correlation between vertices dependent on age and sex (not BMI) => LMM always result in loss of power



Conclusions

LMM offer a more parsimonious and conservative characterisation of the brain-trait associations
Yet the inflated FWER suggests replication may be necessary to safely conclude about an association

OSCA

OmicS-data-based Complex trait Analysis

GCTA SMR GSMD OSCA

GCTB Program in PCTG CTG forum

About
Credits
Questions and Help Requests
Citation
Download

About

OSCA (OmicS-data-based Complex trait Analysis) is a software tool written in C/C++ for the analysis of complex traits using multi-omics data. It is developed by [Futao Zhang](#) and [Jian Yang](#) at [Institute for Molecular Bioscience](#), The University of Queensland. Bug reports or questions: Jian Yang <jian.yang@uq.edu.au> or Futao Zhang <futao.zhang@imb.uq.edu.au>.

<https://cnsgenomics.com/software/osca/#Overview>

For all analyses, including data management,
phenotype simulation and LMM from any OMICs,
brain imaging or big-data



Futao Zhang



Jian Yang

Thank you

To all participants and data managers of the UKB and OASIS studies
To the ISBI 2020 e-participants and organising committee

The University of Queensland

Peter Visscher

Naomi Wray

Jian Yang

Futao Zhang

Kathryn E. Kemper

Julia Sidorenko

OSCA software

<https://cnsgenomics.com/software/osca/#Overview>

Paris Brain Institute

Olivier Colliot

Contact: b.couvyduchesne@uq.edu.au