



PCA outperforms popular hidden variable inference methods for molecular QTL mapping

Speaker: LI Yuekai

Date: 07/11/2025





PCA outperforms popular hidden variable inference methods for molecular QTL mapping

Speaker: LI Yuekai

Date: 07/11/2025



CONTENTS

1. Introduction

2. Methods

3. Results

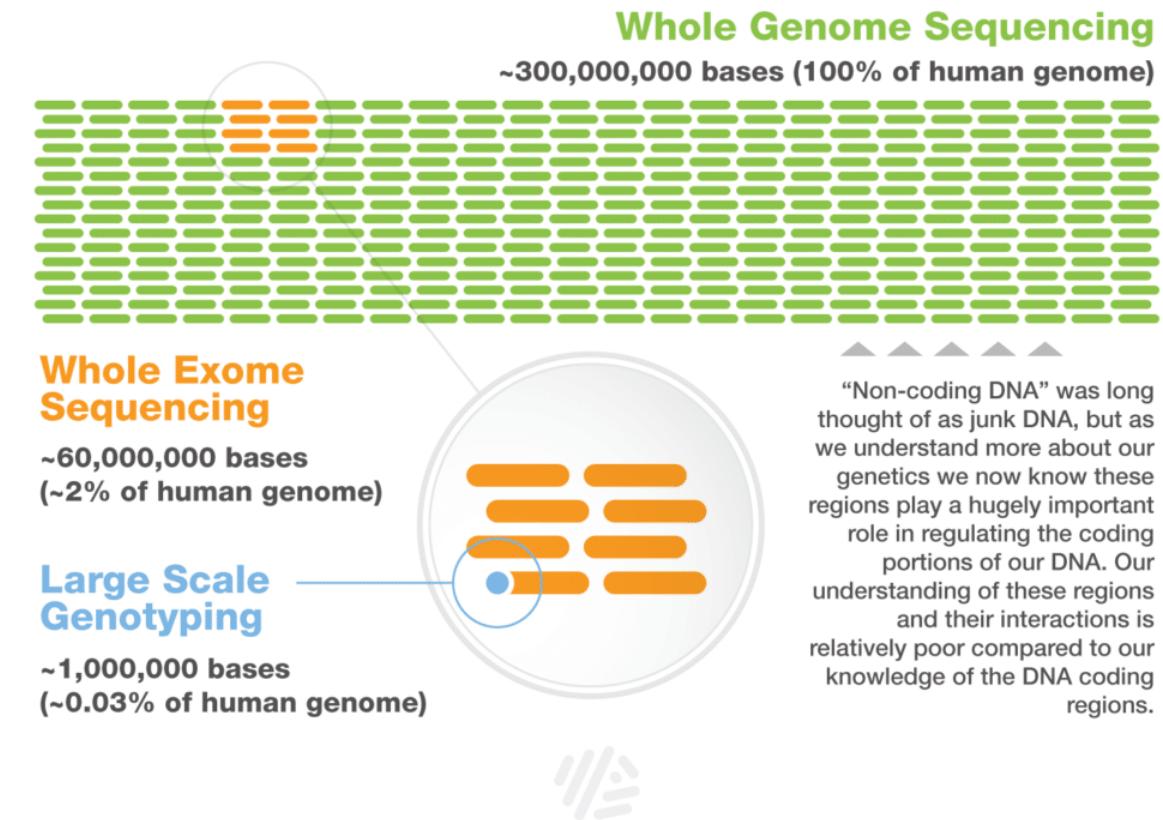
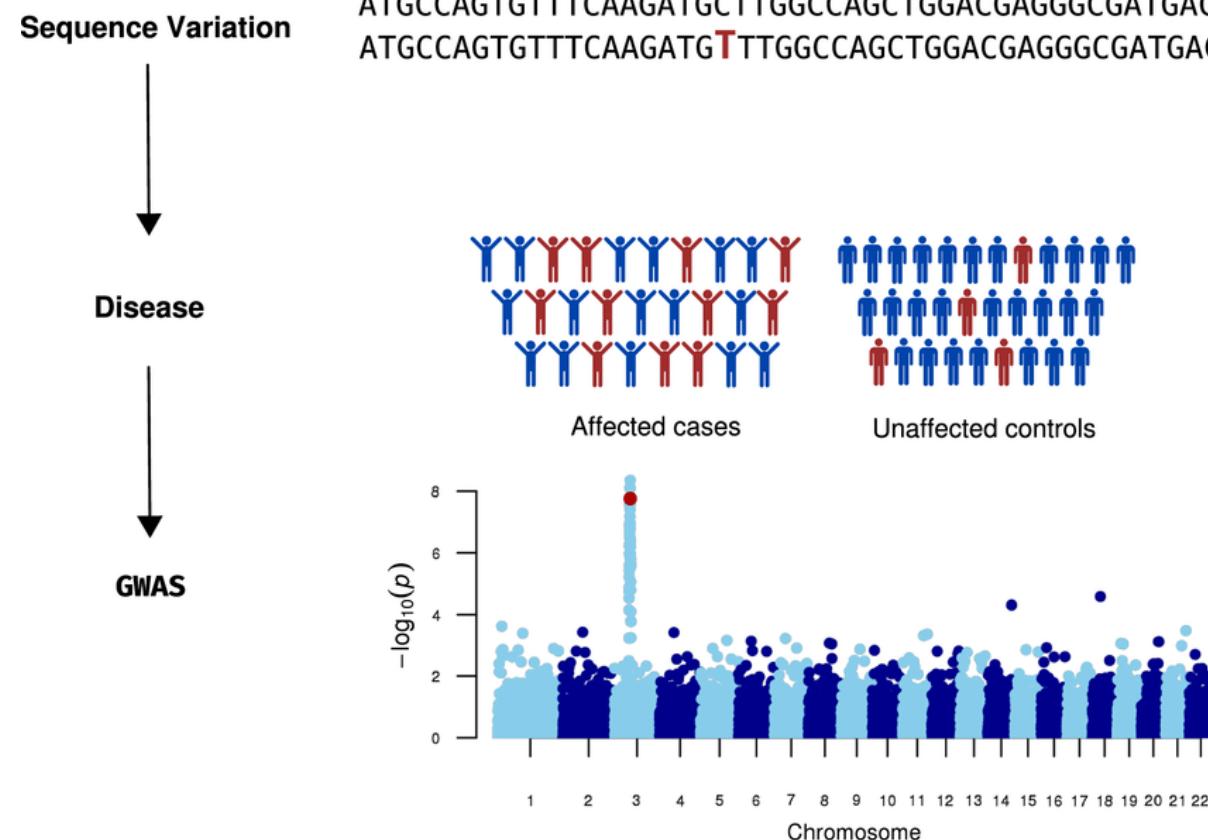
4. Discussion

01

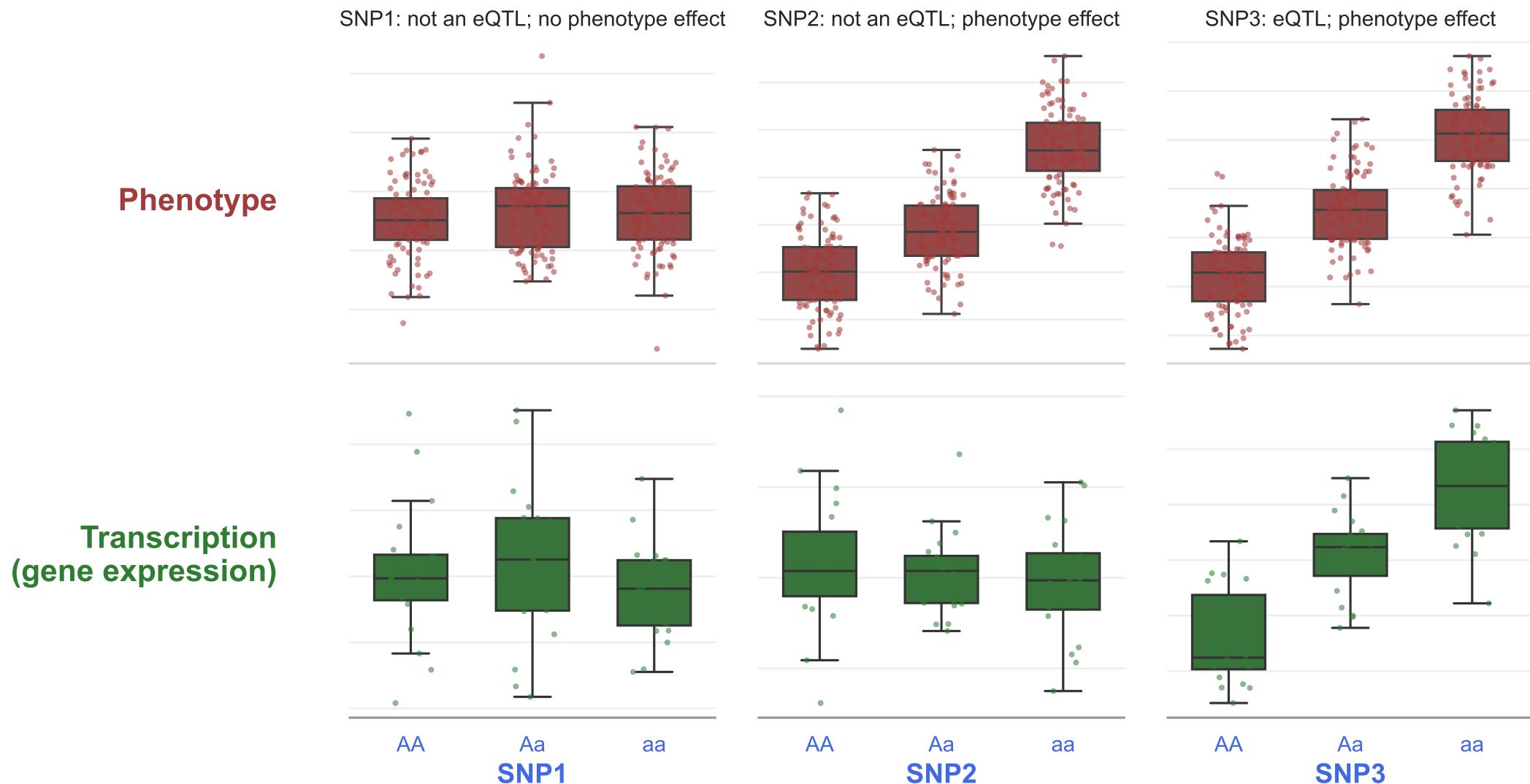
Introduction



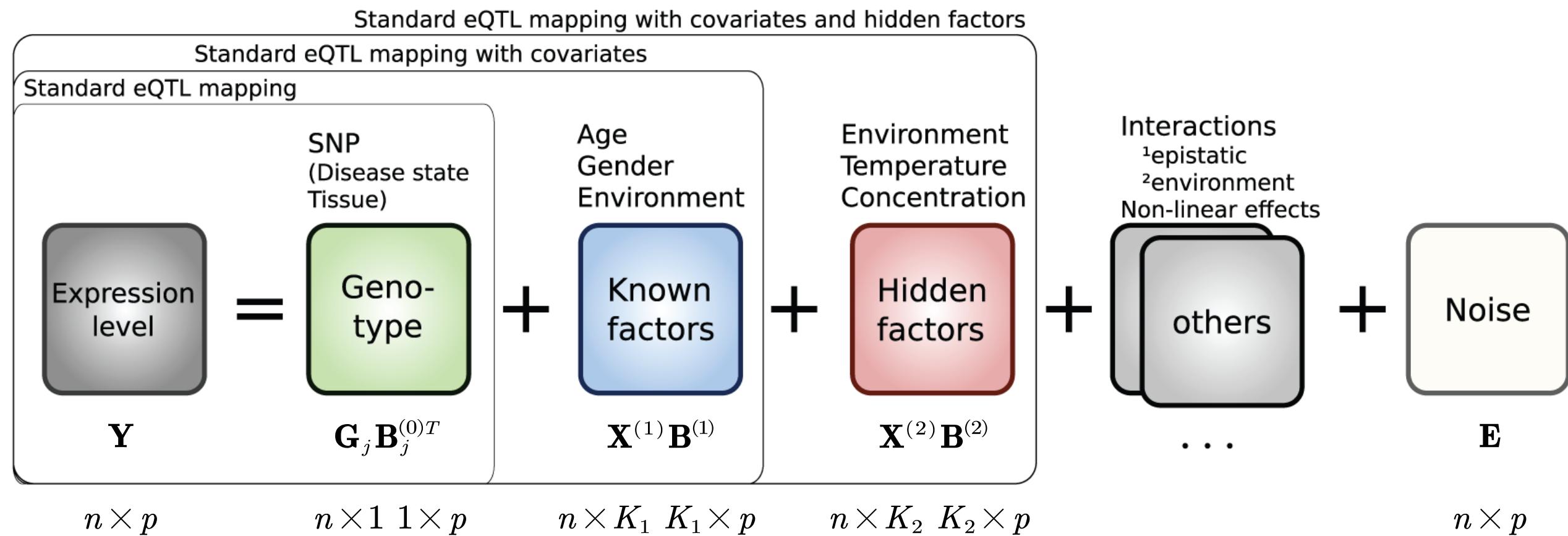
■ GWAS and complex regulatory mechanisms of SNPs



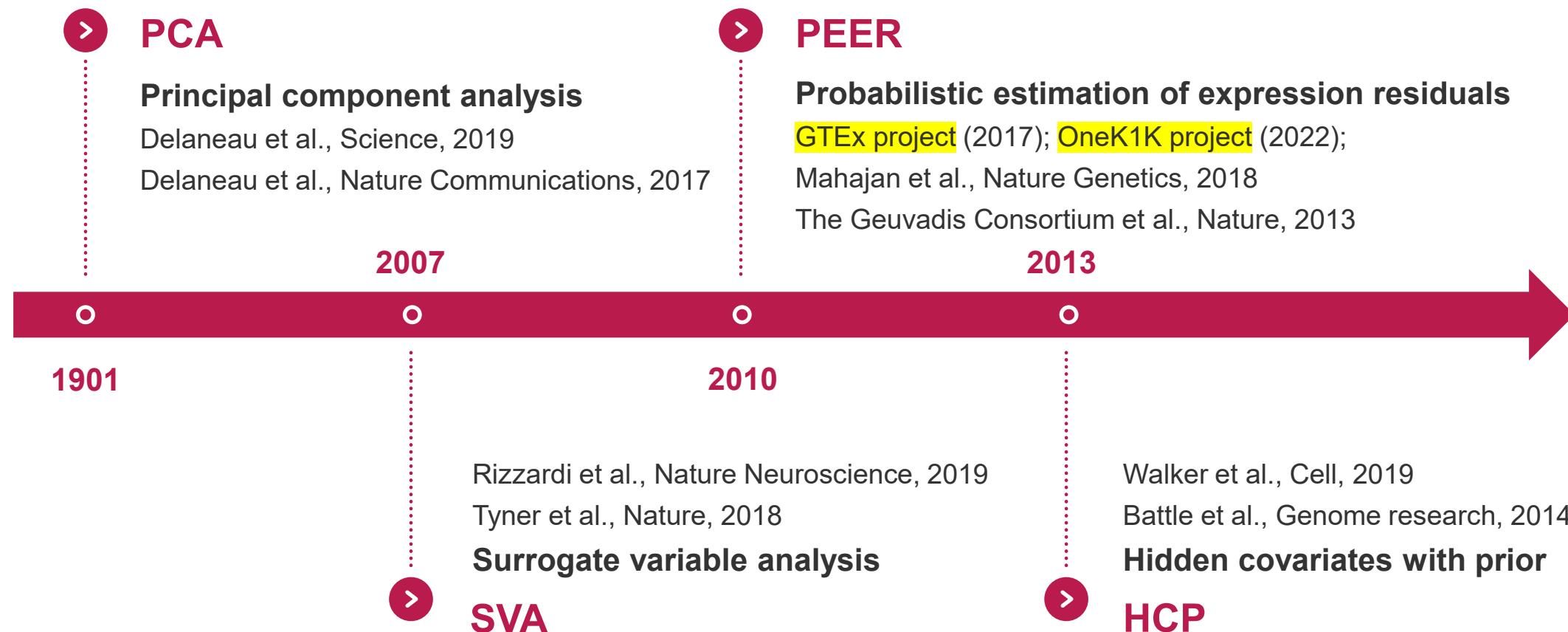
The expression quantitative trait loci (eQTL) analysis



■ The standard eQTL mapping model



■ Some popular hidden variable inference methods for eQTL mapping



02

Methods



■ Principal component analysis (PCA)

Basic idea: Find a loading matrix \mathbf{W} , which linearly transforms the observed data \mathbf{Y} into a new feature space.

$$\mathbf{X} = \mathbf{Y} \mathbf{W}$$

$$n \times p \quad n \times p \quad p \times p$$

Constraints: Constrain \mathbf{W} to be orthogonal, and choose its columns \mathbf{W}_j so that each transformed coordinate $\mathbf{Y}\mathbf{W}_j$ sequentially has the maximum variance.

$$\mathbf{W}_1 = \underset{\mathbf{W}_1^* \in \mathbb{R}^p}{\operatorname{argmax}} \operatorname{Var}(\mathbf{Y}\mathbf{W}_1^*) \quad \text{subject to } \|\mathbf{W}_1^*\|_2 = 1,$$

$$\mathbf{W}_2 = \underset{\mathbf{W}_2^* \in \mathbb{R}^p}{\operatorname{argmax}} \operatorname{Var}(\mathbf{Y}\mathbf{W}_2^*) \quad \text{subject to } \|\mathbf{W}_2^*\|_2 = 1, \mathbf{W}_2^{*T} \mathbf{W}_1 = 0,$$

$$\vdots$$

$$\mathbf{W}_p = \underset{\mathbf{W}_p^* \in \mathbb{R}^p}{\operatorname{argmax}} \operatorname{Var}(\mathbf{Y}\mathbf{W}_p^*) \quad \text{subject to } \|\mathbf{W}_p^*\|_2 = 1, \mathbf{W}_p^{*T} \mathbf{W}_j = 0, \forall j < p.$$

Result: Keep only the first K columns ($K < p$) of \mathbf{W} to obtain a low-dimensional representation.

$$\tilde{\mathbf{X}} = \mathbf{Y} \tilde{\mathbf{W}}$$

$$n \times K \quad n \times p \quad p \times K$$

■ Principal component analysis (PCA)

Algorithm 1 Principal component analysis (PCA) algorithm

- 1: **Input:** Data matrix $\mathbf{Y} \in \mathbb{R}^{n \times p}$ (already standardized); cumulative variance threshold $\tau \in (0, 1)$
- 2: **Output:** Low-dimension representation $\tilde{\mathbf{X}} \in \mathbb{R}^{n \times K}$; number of components K
- 3: Compute sample covariance matrix: $\mathbf{S} \leftarrow \frac{1}{n-1} \mathbf{Y}^\top \mathbf{Y}$
- 4: Eigendecomposition: $\mathbf{S} = \mathbf{V} \Lambda \mathbf{V}^\top$ $\triangleright \Lambda = \text{diag}(\lambda_1, \dots, \lambda_p), \lambda_1 \geq \dots \geq \lambda_p$
- 5: Compute cumulative variance ratio: $r_k \leftarrow \sum_{j=1}^k \lambda_j / \sum_{j=1}^p \lambda_j, k = 1, \dots, p$
- 6: Determine smallest K such that $r_K \geq \tau$
- 7: Select leading eigenvectors: $\tilde{\mathbf{W}} \leftarrow \mathbf{V}_{(:,1:K)}$
- 8: Project the observed data matrix: $\tilde{\mathbf{X}} \leftarrow \mathbf{Y} \tilde{\mathbf{W}}$
- 9: **return** $\tilde{\mathbf{X}}, K$

$$\sum_{j=1}^p \text{Var}(\mathbf{Y}_j) = \sum_{j=1}^p \lambda_j = \sum_{j=1}^p \text{Var}(\mathbf{X}_j)$$

■ Factor analysis (FA)

Basic model:

$$\mathbf{Y} = \mathbf{X} \mathbf{W} + \mathbf{E}$$

$n \times p \quad n \times K \quad K \times p \quad n \times p$

Normality assumption:

Hidden factor matrix:	$x_{ik} \stackrel{\text{iid.}}{\sim} \mathcal{N}(0, 1), \quad i = 1, \dots, n; \quad k = 1, \dots, K$]
Error matrix:	$e_{ij} \stackrel{\text{ind.}}{\sim} \mathcal{N}(0, \sigma_j^2), \quad i = 1, \dots, n; \quad j = 1, \dots, p$	

$\text{Cov}(x_{ik}, e_{ij}) = 0, \text{ for } \forall i, j, k$

where:

\mathbf{Y} is the observed data, and the marginal distribution of \mathbf{Y} is normal.

\mathbf{W} and σ_j^2 are fixed but unknown parameters

\mathbf{X} is the missing data (latent variable)

} As shown by Rubin and Thayer (1997), the EM algorithm can be used to solve this maximization problem by treating \mathbf{X} as missing data.

■ Factor analysis (FA)

Joint distribution:

Define $\Psi = [\sigma_1^2, \dots, \sigma_p^2]$, we have:

$$\begin{bmatrix} \mathbf{x}_i \\ \mathbf{y}_i \end{bmatrix} \stackrel{\text{iid.}}{\sim} \mathcal{N} \left(\mathbf{0}, \begin{bmatrix} \mathbf{I}_K & \mathbf{W} \\ \mathbf{W}^T & \mathbf{W}^T \mathbf{W} + \Psi \end{bmatrix} \right), \quad i = 1, \dots, n$$

EM algorithm:

E-step:

$$\mathbb{E}[\mathbf{x}_i | \mathbf{y}_i] = \mathbf{W}(\mathbf{W}^T \mathbf{W} + \Psi)^{-1} \mathbf{y}_i$$

$$\text{Cov}[\mathbf{x}_i | \mathbf{y}_i] = \mathbf{I}_K - \mathbf{W}(\mathbf{W}^T \mathbf{W} + \Psi)^{-1} \mathbf{W}^T$$

$$\mathbb{E}[\mathbf{x}_i \mathbf{x}_i^T | \mathbf{y}_i] = \text{Cov}[\mathbf{x}_i | \mathbf{y}_i] + \mathbb{E}[\mathbf{x}_i | \mathbf{y}_i] \mathbb{E}[\mathbf{x}_i | \mathbf{y}_i]^T$$

Thomson's factor scores



M-step:

$$\mathbf{W}^{\text{new}} = \left(\frac{1}{n} \sum_{i=1}^n \mathbb{E}[\mathbf{x}_i | \mathbf{y}_i] \mathbf{y}_i^T \right) \left(\frac{1}{n} \sum_{i=1}^n \mathbb{E}[\mathbf{x}_i \mathbf{x}_i^T | \mathbf{y}_i] \right)^{-1}$$

$$\Psi^{\text{new}} = \text{diag} \left\{ \frac{1}{n} \sum_{i=1}^n [\mathbf{y}_i \mathbf{y}_i^T - \mathbf{W}^{\text{new} T} \mathbb{E}[\mathbf{x}_i | \mathbf{y}_i] \mathbf{y}_i^T] \right\}$$

parameter estimates at convergence

■ Probabilistic principal component analysis (PPCA)

Basic model:

$$\mathbf{Y} = \mathbf{X} \mathbf{W} + \mathbf{E}$$

$n \times p \quad n \times K \quad K \times p \quad n \times p$

Normality assumption:

$$\left. \begin{array}{l} \text{Hidden factor matrix: } x_{ik} \stackrel{\text{iid.}}{\sim} \mathcal{N}(0, 1), \quad i = 1, \dots, n; \quad k = 1, \dots, K \\ \text{Error matrix: } e_{ij} \stackrel{\text{iid.}}{\sim} \mathcal{N}(0, \sigma^2), \quad i = 1, \dots, n; \quad j = 1, \dots, p \end{array} \right\} \text{Cov}(x_{ik}, e_{ij}) = 0, \text{ for } \forall i, j, k$$

In contrast to factor analysis, the maximum likelihood can be obtained analytically in PPCA:

$$\left. \begin{array}{l} \sigma_{\text{MLE}}^2 = \frac{1}{p - K} \sum_{j=K+1}^p \lambda_j \\ \mathbf{W}_{\text{MLE}} = \mathbf{V}_K (\Lambda_K - \sigma_{\text{MLE}}^2 \mathbf{I}_K)^{\frac{1}{2}} \end{array} \right\} \rightarrow \mathbb{E}[\mathbf{x}_i | \mathbf{y}_i] = \mathbf{W}_{\text{MLE}} (\mathbf{W}_{\text{MLE}}^T \mathbf{W}_{\text{MLE}} + \sigma_{\text{MLE}}^2 \mathbf{I}_p)^{-1} \mathbf{y}_i$$

where, \mathbf{V}_k , Λ_K , λ_j are all derived from the eigendecomposition of the sample covariance matrix.

■ Probabilistic estimation of expression residuals (PEER)

Basic model:

$$\mathbf{Y} = [\mathbf{X}^{(1)}, \mathbf{X}^{(2)}] \mathbf{W} + \mathbf{E}$$

$n \times p$ $n \times (K_1 + K_2)$ $(K_1 + K_2) \times p$ $n \times p$

Normality assumption:

Hidden factor matrix:

$$x_{ik}^{(2)} \stackrel{\text{iid.}}{\sim} \mathcal{N}(0, 1), \quad i = 1, \dots, n; k = 1, \dots, K_2$$

Covariate-specific weight matrix:

$$w_{kj} \stackrel{\text{ind.}}{\sim} \mathcal{N}\left(0, \frac{1}{\beta_k}\right), \quad j = 1, \dots, p; k = 1, \dots, K_1 + K_2$$

Error matrix:

$$e_{ij} \stackrel{\text{ind.}}{\sim} \mathcal{N}\left(0, \frac{1}{\tau_j}\right), \quad i = 1, \dots, n; j = 1, \dots, p$$

Prior distribution:

$$\beta_k \stackrel{\text{iid.}}{\sim} \Gamma(a_1, b_1), \quad k = 1, \dots, K_1 + K_2$$

$$\tau_j \stackrel{\text{iid.}}{\sim} \Gamma(a_2, b_2), \quad j = 1, \dots, p$$

■ Hidden covariates with prior (HCP)

Basic model:

$$\mathbf{Y} = \mathbf{X}^{(1)} \mathbf{W}^{(1)} + \mathbf{X}^{(2)} \mathbf{W}^{(2)} + \mathbf{E}$$

$n \times p$ $n \times K_1$ $K_1 \times p$ $n \times K_2$ $K_2 \times p$ $n \times p$

- $\mathbf{X}^{(1)}$ is known covariates;
- $\mathbf{X}^{(2)}$ is unknown covariates;
- $\mathbf{W}^{(1)}$ is **fixed effects** for known covariates.

Normality assumption:

$$x_{ik_2}^{(2)} \stackrel{\text{iid.}}{\sim} \mathcal{N}\left((\mathbf{X}^{(1)} \mathbf{C})_{ik_2}, \frac{1}{\lambda_1}\right), \quad i=1, \dots, n; \quad k_2=1, \dots, K_2$$

$$c_{k_1 k_2} \stackrel{\text{iid.}}{\sim} \mathcal{N}\left(0, \frac{1}{\lambda_2}\right), \quad k_1=1, \dots, K_1; \quad k_2=1, \dots, K_2$$

$$w_{k_2 j}^{(2)} \stackrel{\text{iid.}}{\sim} \mathcal{N}\left(0, \frac{1}{\lambda_3}\right), \quad j=1, \dots, p; \quad k_2=1, \dots, K_2$$

$$e_{ij} \stackrel{\text{iid.}}{\sim} \mathcal{N}\left(0, \frac{1}{\lambda}\right), \quad i=1, \dots, n; \quad j=1, \dots, p$$

The prior of unknown covariates depends on linear combinations of known covariates.

Penalized regression form:

$$\operatorname{argmin}_{\mathbf{X}^{(2)}, \mathbf{C}, \mathbf{W}^{(1)}, \mathbf{W}^{(2)}} \|\mathbf{Y} - \mathbf{X}^{(1)} \mathbf{W}^{(1)} - \mathbf{X}^{(2)} \mathbf{W}^{(2)}\|_2^2 + \lambda_1 \|\mathbf{X}^{(2)} - \mathbf{X}^{(1)} \mathbf{C}\|_2^2 + \lambda_2 \|\mathbf{C}\|_2^2 + \lambda_3 \|\mathbf{W}^{(2)}\|_2^2$$

■ Surrogate variable analysis (SVA)

Algorithm 2-1 Iteratively reweighted surrogate variable analysis (IRW-SVA)

```

1: Input:
    • Gene expression matrix  $\mathbf{Y}$ ; variables of interest  $\mathbf{G}$ ; known covariates  $\mathbf{X}_1$  (optional).
    • Number of inferred variables  $K$ ; number of iterations  $B$ .
2: Output:  $\mathbf{X}_2$  ( $n \times K$  matrix of inferred covariates).
3: for each gene do
4:     Regress the  $j$ -th column of  $\mathbf{Y}$  against  $(\mathbf{1}, \mathbf{G}, \mathbf{X}_1)$ ;
5:     Replace the  $j$ -th column of  $\mathbf{Y}$  with the residuals to get  $\mathbf{R}$ ;
6: end for
7: Perform PCA on  $\mathbf{R}$  to get the initial PCs;
8: for  $b \leftarrow 1$  to  $B$  do
9:     for each gene do
10:        Test  $H_0$ : coefficients of  $\mathbf{G} = 0$  in the regression against  $(\mathbf{1}, \mathbf{G}, \mathbf{X}_1, \text{PCs})$  via partial  $F$ -test;
11:        Convert the  $p$ -values to local false discovery rates, denoted as  $\text{lfdr}_1$ ;
12:    end for
13:    for each gene do
14:        Test  $H_0$ : coefficients of PCs = 0 in the regression against  $(\mathbf{1}, \mathbf{X}_1, \text{PCs})$  via partial  $F$ -test;
15:        Convert the  $p$ -values to local false discovery rates, denoted as  $\text{lfdr}_2$ ;
16:    end for
17:    Weight the columns of  $\mathbf{Y}$  by  $\text{lfdr}_1 \times (1 - \text{lfdr}_2)$ ;
18:    Perform PCA on the weighted  $\mathbf{Y}$  after centering;
19: end for
20: return the first  $K$  PCs from the last PCA.

```

SVA is purely algorithmic, based on the PCA algorithm.



■ Surrogate variable analysis (SVA)

Algorithm 2-1 Iteratively reweighted surrogate variable analysis (IRW-SVA)

```

1: Input:
    • Gene expression matrix  $\mathbf{Y}$ ; variables of interest  $\mathbf{G}$ ; known covariates  $\mathbf{X}_1$  (optional).
    • Number of inferred variables  $K$ ; number of iterations  $B$ .
2: Output:  $\mathbf{X}_2$  ( $n \times K$  matrix of inferred covariates).
3: for each gene do
4:     Regress the  $j$ -th column of  $\mathbf{Y}$  against  $(\mathbf{1}, \mathbf{G}, \mathbf{X}_1)$ ;
5:     Replace the  $j$ -th column of  $\mathbf{Y}$  with the residuals to get  $\mathbf{R}$ ;
6: end for
7: Perform PCA on  $\mathbf{R}$  to get the initial PCs;
8: for  $b \leftarrow 1$  to  $B$  do
9:     for each gene do
10:        Test  $H_0$ : coefficients of  $\mathbf{G} = 0$  in the regression against  $(\mathbf{1}, \mathbf{G}, \mathbf{X}_1, \text{PCs})$  via partial  $F$ -test;
11:        Convert the  $p$ -values to local false discovery rates, denoted as  $\text{lfdr}_1$ ;
12:    end for
13:    for each gene do
14:        Test  $H_0$ : coefficients of PCs = 0 in the regression against  $(\mathbf{1}, \mathbf{X}_1, \text{PCs})$  via partial  $F$ -test;
15:        Convert the  $p$ -values to local false discovery rates, denoted as  $\text{lfdr}_2$ ;
16:    end for
17:    Weight the columns of  $\mathbf{Y}$  by  $\text{lfdr}_1 \times (1 - \text{lfdr}_2)$ ;
18:    Perform PCA on the weighted  $\mathbf{Y}$  after centering;
19: end for
20: return the first  $K$  PCs from the last PCA.

```

SVA is purely algorithmic, based on the PCA algorithm.



Algorithm 2-2 The calculation for lfdr in IRW-SVA algorithm

```

1: Input:  $p$ -values  $\text{pVals}$ ; vector of length  $p$ ; threshold  $\lambda = 0.8, \epsilon = 10^{-8}$ .
2: Output: local false discovery rates  $\text{lfdrs}$ .
3:  $\hat{\pi}_0 \leftarrow \text{sum}(\text{pVals} > \lambda)/p(1 - \lambda)$ ;
4: Floor each element of  $\text{pVals}$  at  $\epsilon$  and cap each element of  $\text{pVals}$  at  $1 - \epsilon$ ;
5:  $\hat{f}_0 \leftarrow \text{dnorm}(\text{qnorm}(\text{pVals}))$ ;
6:  $\hat{f}$  is obtained by fitting a smooth curve to the histogram of  $\text{pVals}$  via kernel density estimation;
7:  $\text{lfdrs} \leftarrow \hat{\pi}_0 \hat{f}_0 / \hat{f}$ ;
8: return  $\text{lfdrs}$ ;

```

Algorithm 2-3 the BE algorithm for choosing K in SVA

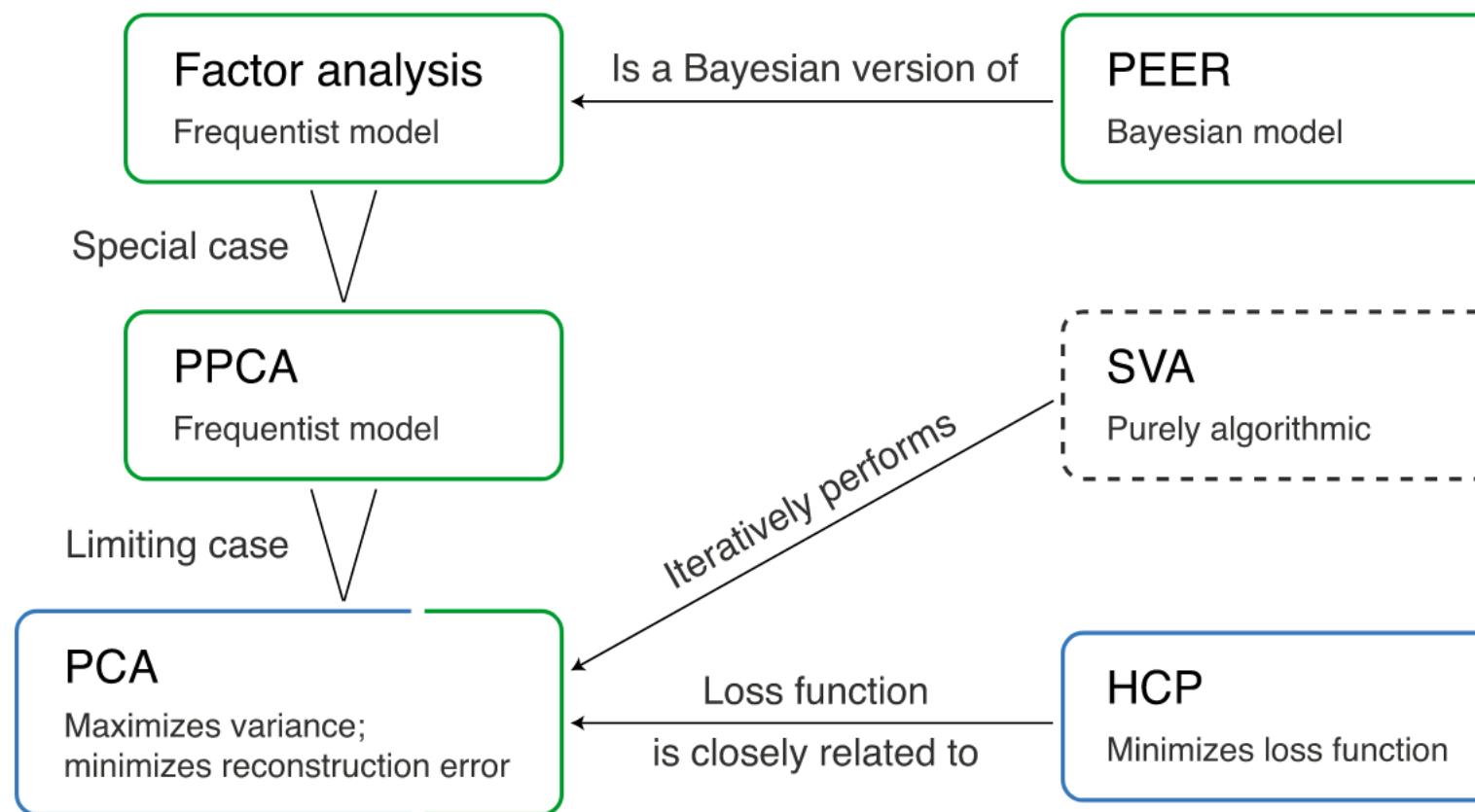
```

1: Input: Gene expression matrix  $\mathbf{Y}$ ; variables of interest  $\mathbf{G}$ ; known covariates  $\mathbf{X}_1$  (optional).
2: Output: Number of hidden covariates  $K$ .
3: Residualize  $\mathbf{Y}$  against  $(\mathbf{1}, \mathbf{G}, \mathbf{X}_1)$  to obtain  $\mathbf{R}$ ;
4: Perform PCA on  $\mathbf{R}$ ;
5: Denote the proportion of variance explained (PVE) by the  $k$ -th PC by  $\text{PVE}_k$ ;
6: for  $b \leftarrow 1$  to  $B = 20$  do
7:     Permute each column of  $\mathbf{R}$  to obtain  $\mathbf{R}_b$ ;
8:     Residualize  $\mathbf{R}_b$  against  $(\mathbf{1}, \mathbf{G}, \mathbf{X}_1)$  to obtain  $\mathbf{R}'_b$ ;
9:     Perform PCA on  $\mathbf{R}'_b$ ;
10:    end for
11:    The  $p$ -value for the  $k$ -th PC is calculated as the proportion of permutations where the PVE of
        the  $k$ -th PC is greater than or equal to  $\text{PVE}_k$ ;
12:    Enforce that the  $p$ -values increase (i.e., are non-decreasing) as  $k$  increases;
13:    return the number of PCs with a  $p$ -value smaller than or equal to  $\alpha = 0.1$ .

```

■ The relationship between PCA and hidden variable inference methods

Classic methods



■ Comparison of PCA, PEER, HCP, and SVA

Method	Variables of interest	Known covariates	K	tuning parameters	probabilistic model
PCA	Unrequired	Unrequired	Adaptive	Unrequired	No
PEER	Unrequired	Optional	Required	Optional	Yes
HCP	Unrequired	Required	Required	Required	Yes
SVA	Required	Optional	Optional	Unrequired	No

03

Results



■ Limitations of the simulation in the original PEER publication

Data analysis limitations:

- (a) The study only compares PEER against the other methods in terms of power, not in terms of FPR or FDR.
- (b) The study does not use PCA or SVA properly.
- (c) The study does not evaluate the different ways of using PEER.
- (d) The study uses non-standard priors for PEER that are different from the default priors.

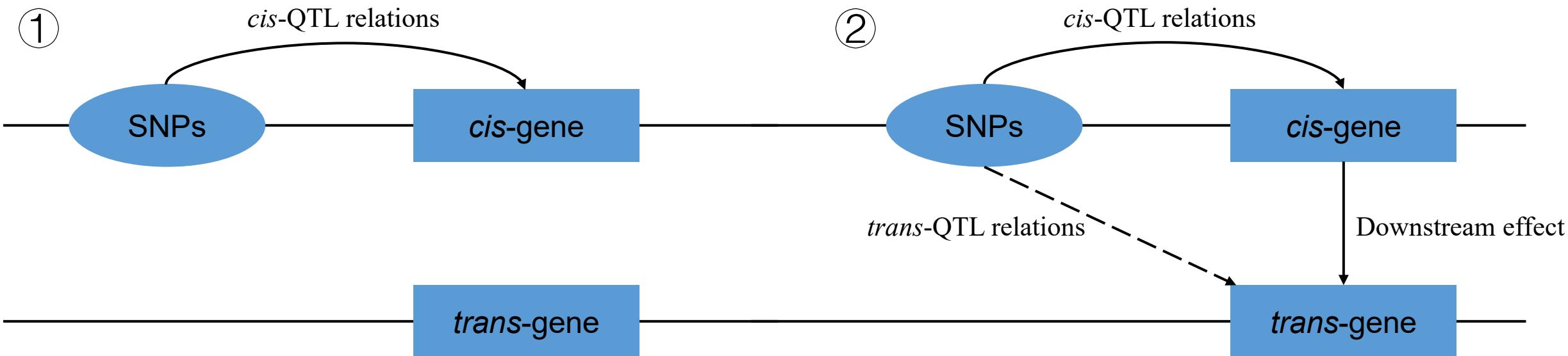
Data simulation limitations:

- (a) The data dimensions are minimal, with only $q = 100$ SNPs in the whole genome.
- (b) The SNP genotypes are simulated independently and identically with a target minor allele frequency (MAF) of 0.4, so there is no linkage disequilibrium (LD) and a higher average MAF than in real data.
- (c) The gene expression levels are primarily driven by *trans*-regulatory effects rather than *cis*-regulatory effects or covariate effects.
- (d) The study only simulates one replicate of one experiment.

■ Simulation design 1 versus Simulation design 2

	Simulation Design 1	Simulation Design 2
Number of experiments	1	176
Number of replicates per experiment	10	2
Genotype data	Simulated (no LD, high MAF)	Real genotype from GTEx
Cis-QTL relations present	✓	✓
Trans-QTL relations present	✓	✗
Number of individuals	n=80	n=838
Number of genes	p=400	p=1000
Number of SNPs	q=100 SNPs in the whole genome	q=1000 local common SNPs per gene
Number of causal SNPs per gene	random	1 or random
Number of known covariates	$k_1=3$	$k_1=2,3,5,$ or 8, respectively
Number of hidden covariates	$k_2=7$	$k_2=3,7,15,$ or 22, respectively

■ Data simulation for trans-QTL detection (design 1)



$$\mathbf{Y} = \mathbf{Y}_{\text{BeforeDSE}} + \mathbf{Y}_{\text{DSE}}$$

$n \times p$ $n \times p$ $n \times p$



$$\mathbf{Y}_{\text{BeforeDSE}} = \mathbf{G} (\mathbf{I} \odot \mathbf{B}) + \mathbf{X}_1 \mathbf{B}_1 + \mathbf{X}_2 \mathbf{B}_2 + \mathbf{E}$$

$n \times p$ $n \times q$ $q \times p$ $q \times p$ $n \times K_1$ $K_1 \times p$ $n \times K_2$ $K_2 \times p$ $n \times p$



$$\mathbf{Y}_{\text{DSE}} = \mathbf{Y}_{\text{BeforeDSE}} (\mathbf{I}_0 \odot \mathbf{B}_0)$$

$n \times p$ $n \times p$ $p \times p$ $p \times p$

■ Data simulation for trans-QTL detection (design 1)

$$\begin{aligned} \mathbf{Y} &= \mathbf{Y}_{\text{BeforeDSE}} + \mathbf{Y}_{\text{DSE}} & \longrightarrow & \mathbf{Y}_{\text{BeforeDSE}} = \mathbf{G} (\mathbf{I} \odot \mathbf{B}) + \mathbf{X}_1 \mathbf{B}_1 + \mathbf{X}_2 \mathbf{B}_2 + \mathbf{E} \\ n \times p & \qquad n \times p & \qquad n \times q & \qquad q \times p & \qquad q \times p & \qquad n \times K_1 & \qquad K_1 \times p & \qquad n \times K_2 & \qquad K_2 \times p & \qquad n \times p \end{aligned}$$

$$\longrightarrow \mathbf{Y}_{\text{DSE}} = \mathbf{Y}_{\text{BeforeDSE}} (\mathbf{I}_0 \odot \mathbf{B}_0)$$

$$n \times p \qquad n \times p \qquad p \times p \qquad p \times p$$

Before Downstream effect (*cis*-QTL):

G: each element is drawn independently from $\text{Bin}(2, 0.4)$.

I: each element is drawn independently from $\text{Ber}(0.01)$.

B: each element is drawn independently from $N(0, 4)$.

X₁, **X**₂: each element is drawn independently from $N(0, 0.6)$.

B₁, **B**₂: $(\mathbf{B}_1)_{kj}, (\mathbf{B}_2)_{kj}$ is drawn from $N(0, \sigma_k^2)$, with the covariate-specific effect size variance $\sigma_k^2 \sim 0.8 * (\Gamma(2.5, 0.6))^2$.

E: E_{ij} is drawn from $N(0, 1/\tau_j)$, with the gene-specific noise precision $\tau_j \sim \Gamma(3, 1)$.

Downstream effect (trans-QTL):

I₀: $p \times p$ binary matrix with 3 rows containing exactly 30 ones (excluding diagonal entries).

B₀: each element is drawn independently from $N(8, 0.8)$ for “strong downstream effects”.

$\mathbf{I}_0 =$	$\begin{bmatrix} 0 & 0 & 0 & \cdots & 0 & 0 & 0 \\ 0 & \ddots & 0 & \cdots & 0 & 0 & 0 \\ 0 & 1 & \ddots & \cdots & 1 & 0 & 1 \\ \vdots & \cdots & \cdots & \mathbf{0} & \cdots & \cdots & \vdots \\ 0 & 0 & 0 & \cdots & \ddots & 0 & 0 \end{bmatrix}$	only 3 rows: sum = 30
	$\begin{bmatrix} 1 & 1 & 0 & \cdots & 0 & \ddots & 1 \\ 1 & 0 & 0 & \cdots & 1 & 1 & 0 \end{bmatrix}$	25

■ Data Simulation for cis-QTL detection (design 2)

$$\mathbf{Y} = \mathbf{G}_{\text{array}} \otimes (\mathbf{I} \odot \mathbf{B}) + \mathbf{X}_1 \mathbf{B}_1 + \mathbf{X}_2 \mathbf{B}_2 + \mathbf{E}$$

$n \times p$ $n \times q \times p$ $q \times p$ $q \times p$ $n \times K_1$ $K_1 \times p$ $n \times K_2$ $K_2 \times p$ $n \times p$

Genotype component:

$\mathbf{G}_{\text{array}}$: $\mathbf{G}_{\text{array}}[, , j]$ is an $n \times q$ real genotype matrix for the q local common SNPs for gene j .

I: each column \mathbf{I}_j is drawn independently from $\text{Mult}(q; 1/q, \dots, 1/q)$, or each element of \mathbf{I}_j is drawn independently from $\text{Ber}(1/q)$.

B: each element is drawn independently from $N(0, 1)$.

Covariate component:

$\mathbf{X}_1, \mathbf{X}_2$: each element is drawn independently from $N(0, 1)$.

$\mathbf{B}_1, \mathbf{B}_2$: each element is drawn independently from $N(0, 1)$ and scaled according to the specified PVE.

Noise component:

\mathbf{E} : each element is drawn independently from $N(0, 1)$ and scaled according to the specified PVE.

■ Evaluation metrics

01

Runtime of the hidden variable inference step

We compare the runtime of 15 methods, including ideal, unadjusted methods and 13 variants of PCA, SVA, PEER, and HCP based on two simulation studies..

02

The area under the PR curve (AUPRC) of the QTL result

We calculate the area under the precision-recall curve (AUPRC). We also compare the powers of the different methods in Simulation Design1.

03

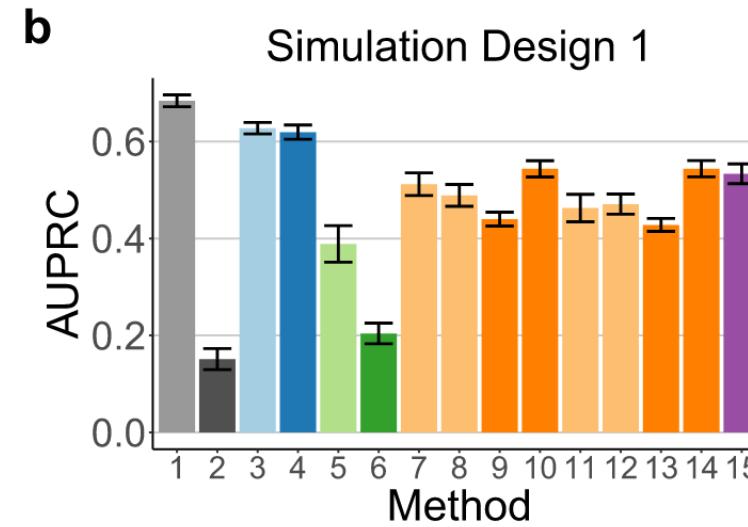
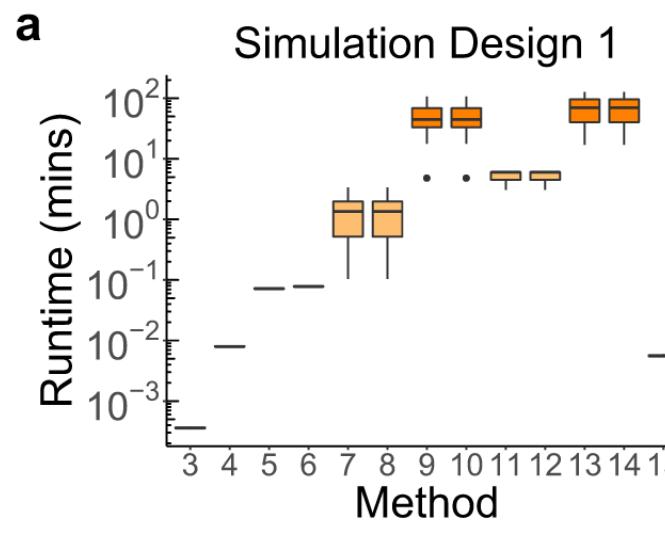
Adjusted R² measures

- adjusted R² score
(regressing each true hidden covariate against the inferred covariates)
 - reverse adjusted R² score
(regressing each inferred covariate against the true hidden covariates)
 - concordance score
(the average of the adjusted R² score and the reverse adjusted R² score)

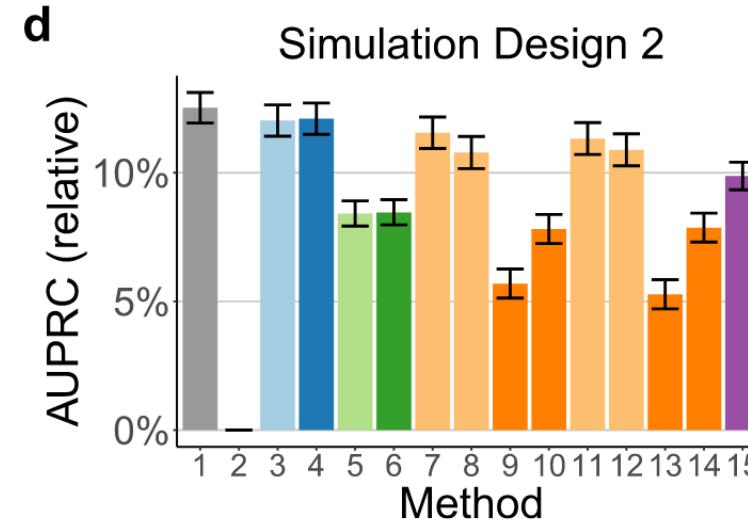
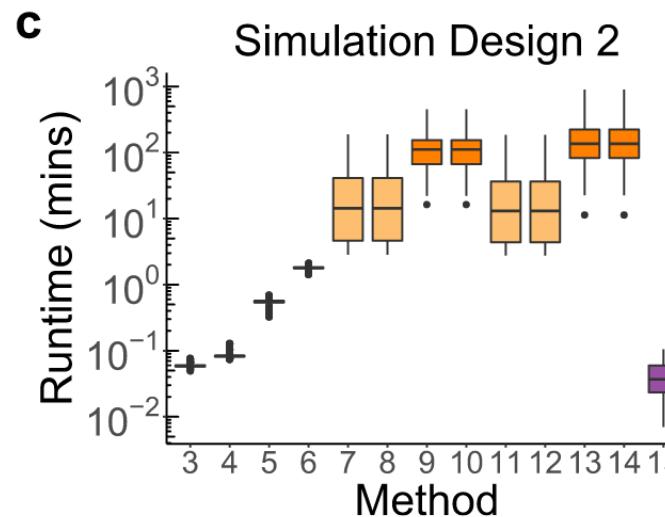
■ Summary of the 15 methods compared simulation studies

	Methods	Response and Covariates	Abbreviations (if selected)
1	Ideal	$\mathbf{Y}, \mathbf{X}_1 + \mathbf{X}_2$	Ideal
2	Unadjusted	\mathbf{Y}, \mathbf{X}_1	Unadjusted
3	PCA _direct_screeK	\mathbf{Y}, \mathbf{X}_1 (filtered) + top PCs	PCA
4	PCA _resid_screeK	$\mathbf{Y}, \mathbf{X}_1 +$ top PCs	
5	SVA _trueK	$\mathbf{Y}, \mathbf{X}_1 +$ SVs	
6	SVA _BE	$\mathbf{Y}, \mathbf{X}_1 +$ SVs	SVA
7	PEER _noCov_trueK_factors	\mathbf{Y}, \mathbf{X}_1 (filtered) + PEER factors	
8	PEER _noCov_trueK_residuals	$\mathbf{Y}_{\text{resid}}, \text{NULL}$	
9	PEER _noCov_largeK_factors	\mathbf{Y}, \mathbf{X}_1 (filtered) + PEER factors	
10	PEER _noCov_largeK_residuals	$\mathbf{Y}_{\text{resid}}, \text{NULL}$	
11	PEER _withCov_trueK_factors	$\mathbf{Y}, \mathbf{X}_1 +$ PEER factors	PEER, true K, factors
12	PEER _withCov_trueK_residuals	$\mathbf{Y}_{\text{resid}}, \text{NULL}$	
13	PEER _withCov_largeK_factors	$\mathbf{Y}, \mathbf{X}_1 +$ PEER factors	
14	PEER _withCov_largeK_residuals	$\mathbf{Y}_{\text{resid}}, \text{NULL}$	PEER, large K, residuals
15	HCP _trueK	$\mathbf{Y}, \mathbf{X}_1 +$ HCPs	HCP

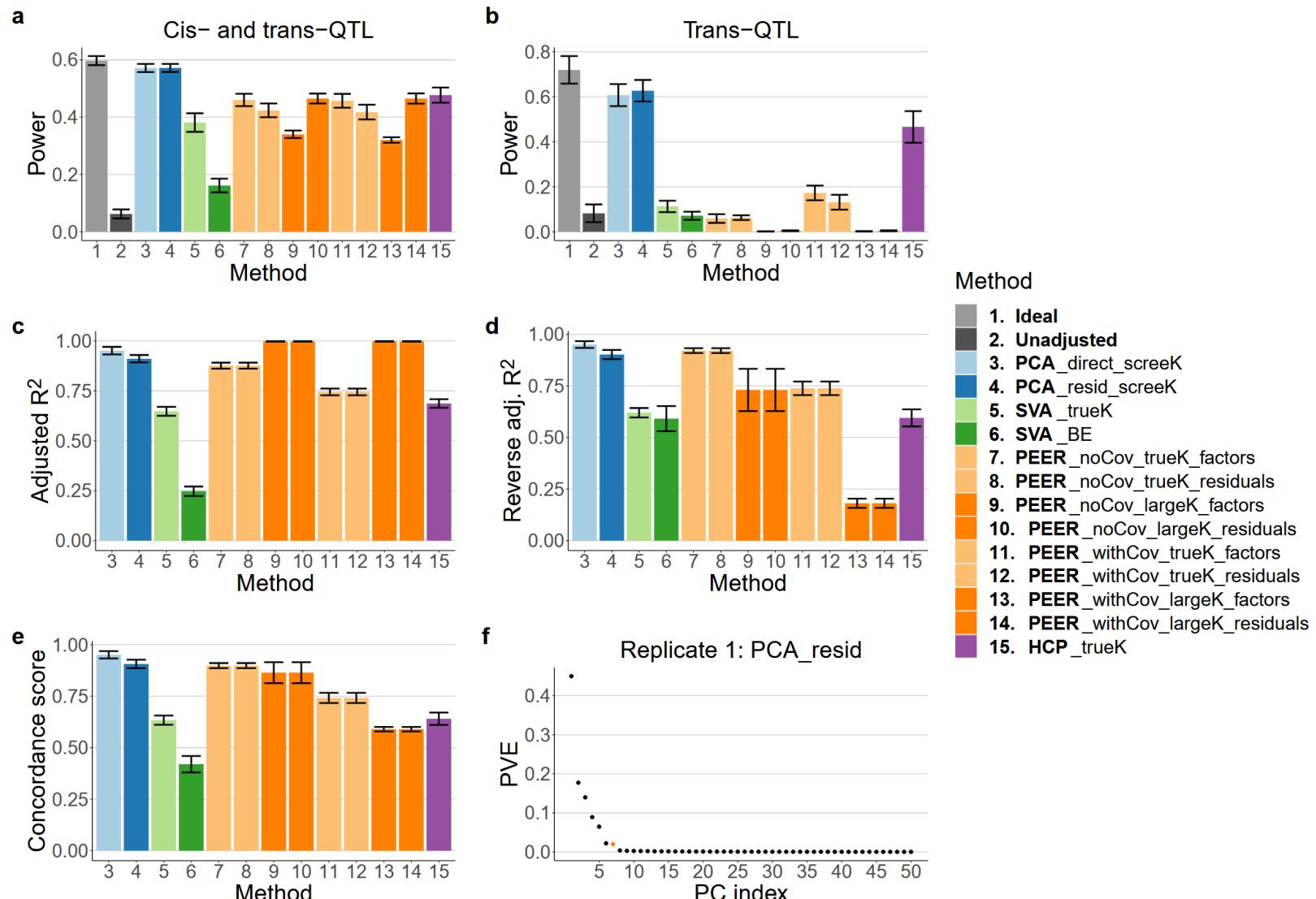
■ Runtime and AUPRC comparison of all 15 methods



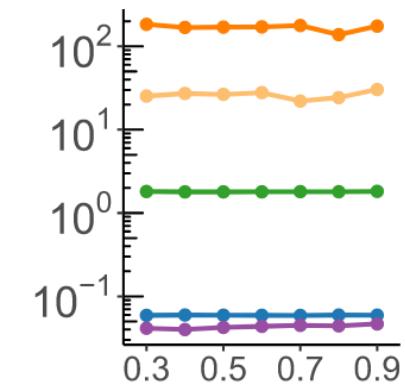
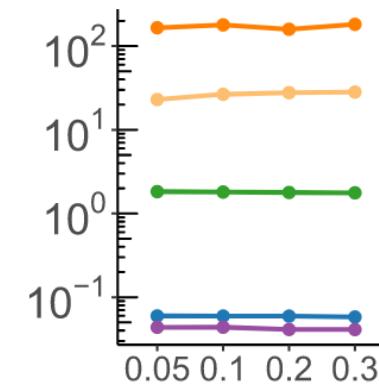
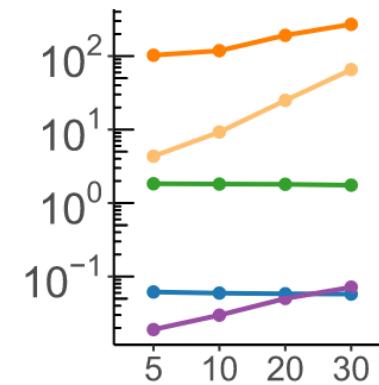
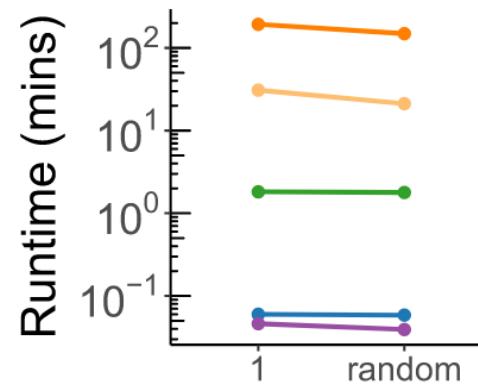
- Method**
1. **Ideal**
 2. **Unadjusted**
 3. **PCA_direct_screeK**
 4. **PCA_resid_screeK**
 5. **SVA_trueK**
 6. **SVA_BE**
 7. **PEER_noCov_trueK_factors**
 8. **PEER_noCov_trueK_residuals**
 9. **PEER_noCov_largeK_factors**
 10. **PEER_noCov_largeK_residuals**
 11. **PEER_withCov_trueK_factors**
 12. **PEER_withCov_trueK_residuals**
 13. **PEER_withCov_largeK_factors**
 14. **PEER_withCov_largeK_residuals**
 15. **HCP_trueK**



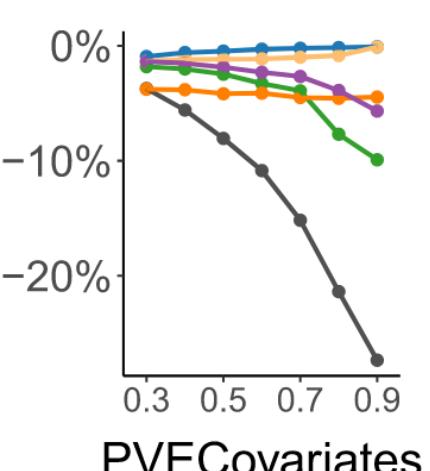
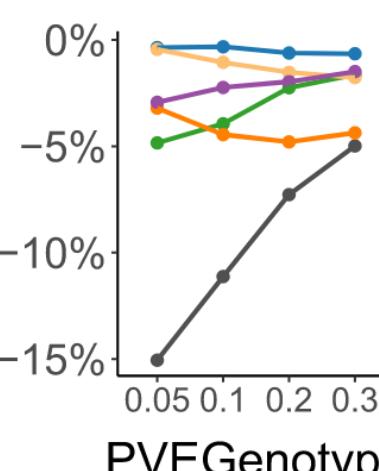
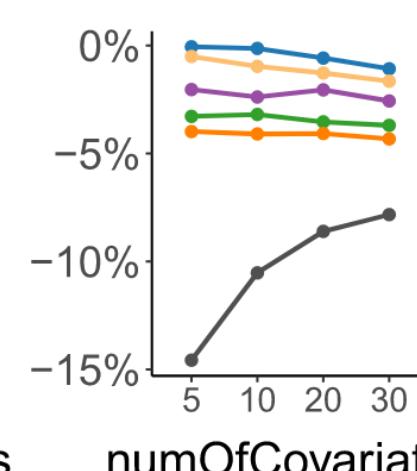
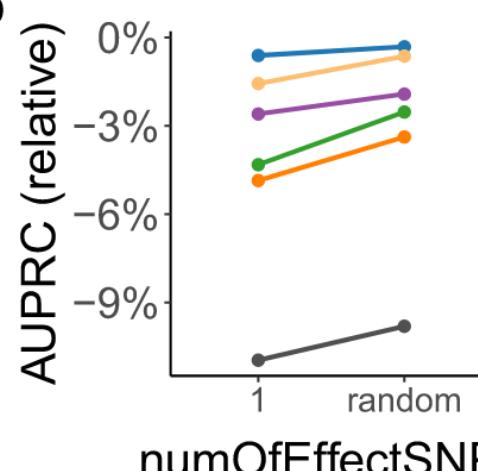
■ Power and adjusted R² measures comparison of all 15 methods in simulation design 1



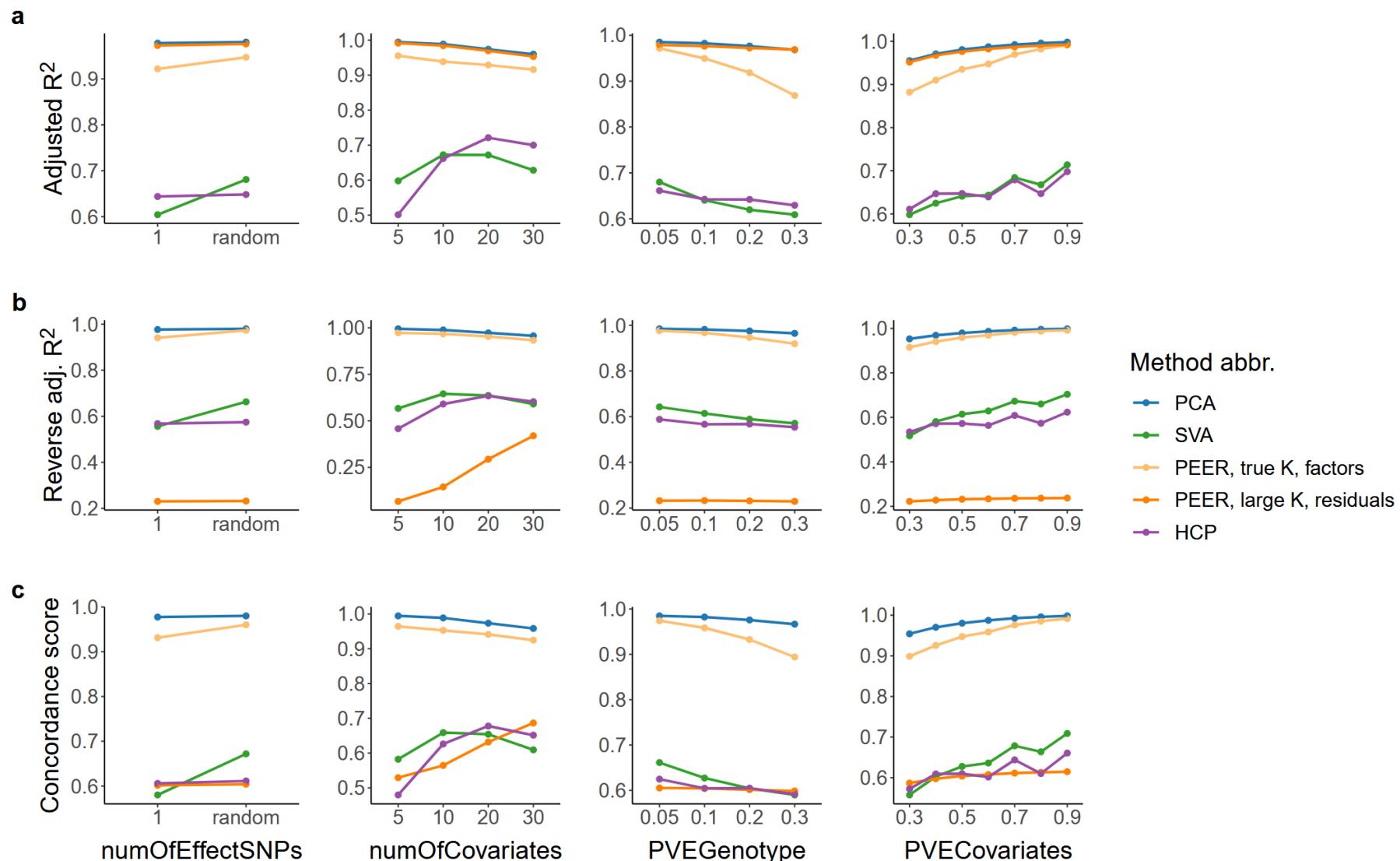
■ Detailed runtime and AUPRC comparison of the selected methods in simulation design 2

a

Method abbr.

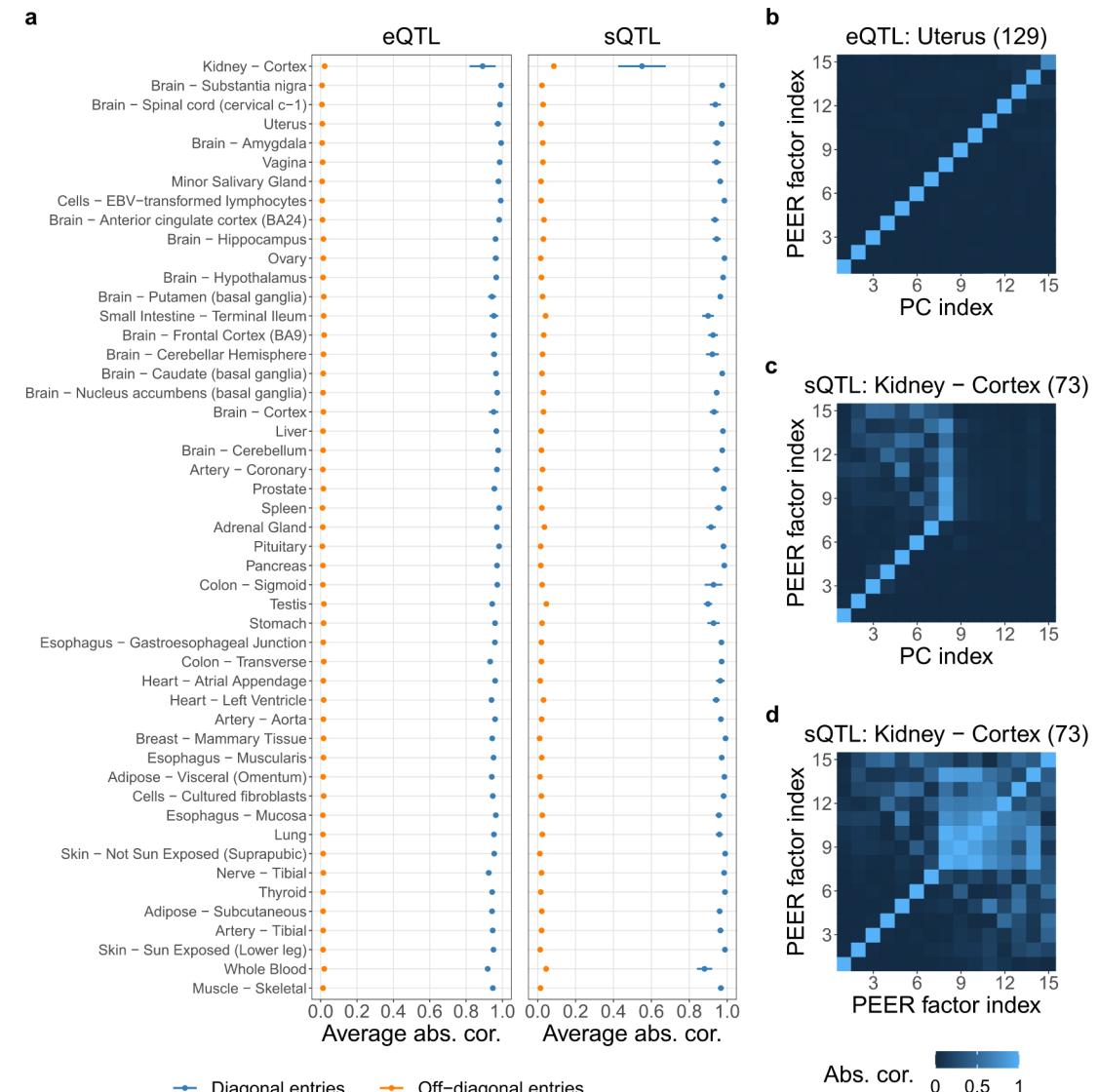
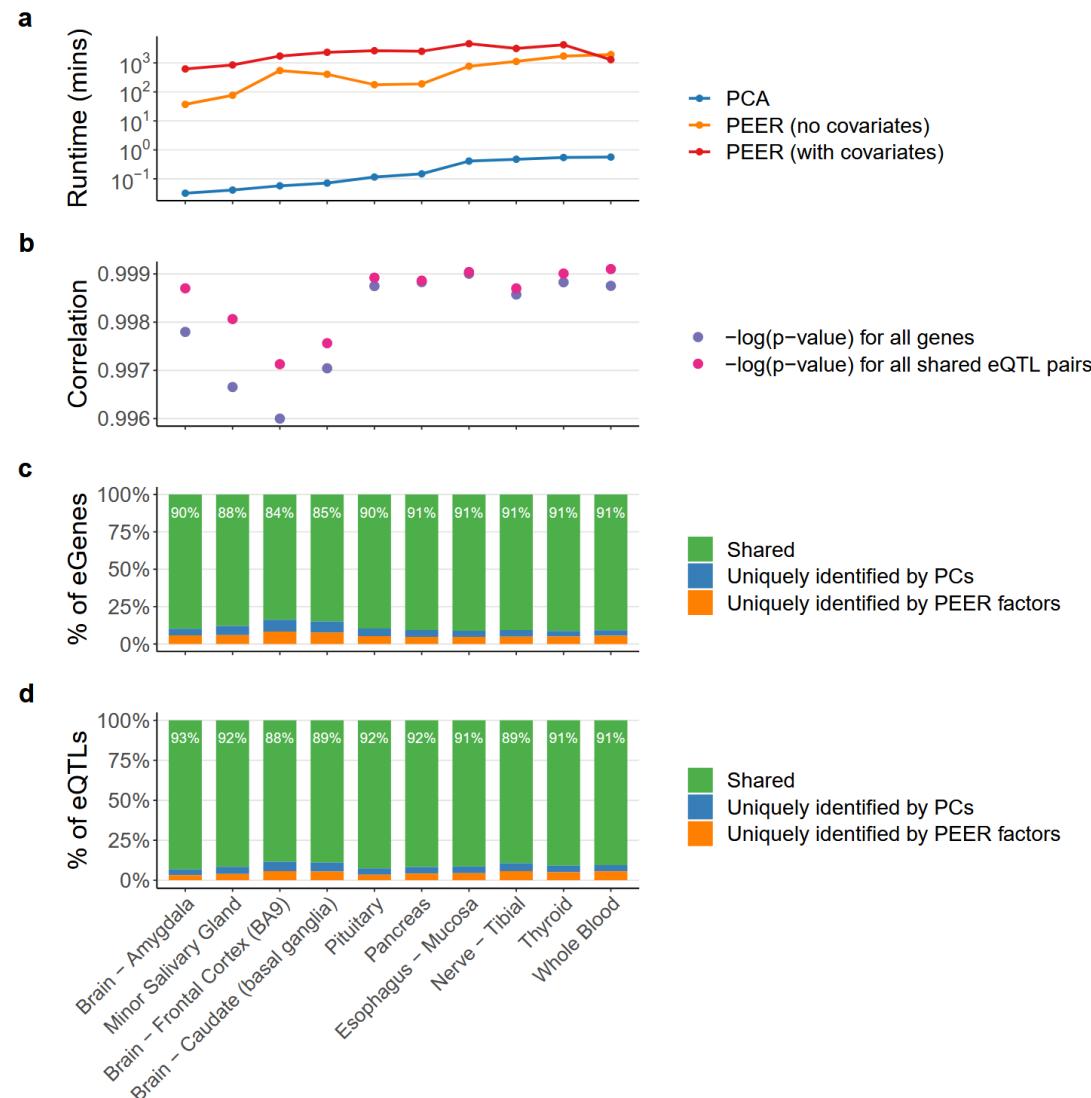
- Unadjusted
- PCA
- SVA
- PEER, true K, factors
- PEER, large K, residuals
- HCP

b


■ Detailed adjusted R² measures comparison of the selected methods in simulation design 2



■ Comparison of PEER factors and PCs in GTEx eQTL and sQTL analysis



■ Comparison of PEER factors and PCs in GTEx eQTL and sQTL analysis

Algorithm S1: Reordering of PEER factors based on PCs (Figure 5).

Inputs:

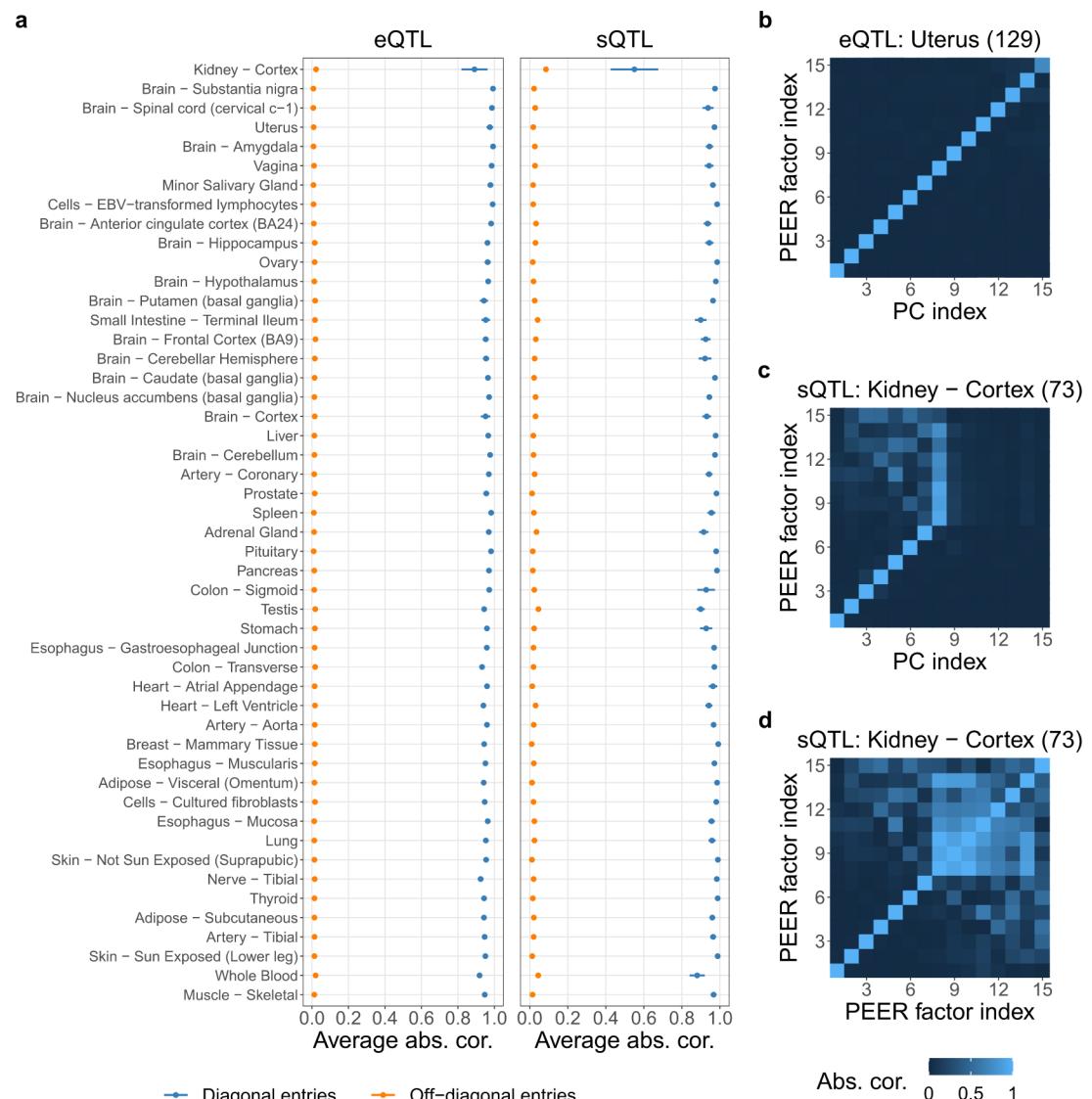
- K PEER factors.
- K PCs.

Output: K PEER factors (reordered).

```

1 for  $k \leftarrow 1$  to  $K$  do
2   | Select the PEER factor that is the most highly correlated with the  $k$ th PC from the PEER
     | factors that have not been selected yet.
3 end
4 return the PEER factors in the order that they were selected in.

```

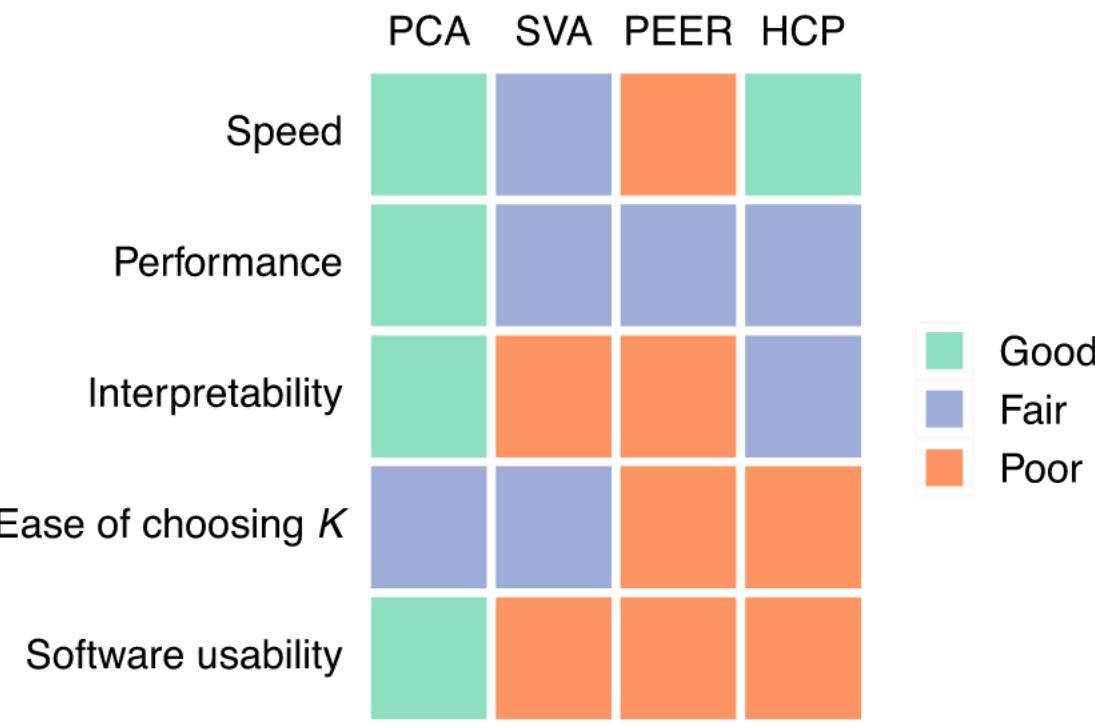


04

Discussion



■ PCA outperforms other methods for the inference of the hidden variables



- PCA is faster.
- PCA is better-performing.
- PCA can be interpreted and used as both a dimension reduction and a factor discovery method.
- PCA offers convenient ways of choosing K.
- There is no consensus on how PEER, HCP and SVA should be used.



香港城市大學
City University of Hong Kong

THANK YOU

Speaker: LI Yuekai

Date: 07/11/2025





香港城市大學
City University of Hong Kong

THANK YOU

Speaker: LI Yuekai

Date: 07/11/2025

