### Modelling molecular heterogeneity between individuals and single cells

### GCC 2015

Oliver Stegle EMBL, European Bioinformatics Institute

### Heterogeneity between individuals and single cells

### variation of interest

confounding





single-cell variation

genetic associations with phenotype

### population variation





differentiation processes

### UK human iPS cell consortium: genotype to phenotype



- Discover how
  - genetic variation affects cellular function
  - genetic lesions lead to disease phenotypes
- In so doing, create an open access resource for the wider biomedical community:
  - HapMap-like consent for ~500 healthy normals
  - Controlled access for ~500 disease samples

### www.hipsci.org

Phil Beales, Ewan Birney, Laura Clarke, Daniel Gaffney, Angus Lamond, Richard Durbin, Fiona Watt



### **Multi-omics association genetics**



### Multi-modal data integration across molecular layers





## Multi-model association genetics: statistical challenges and opportunities

- **Challenge**: Large-scale multiple testing problem:
  - Need to consider potentially millions of loci and adjust for multiple testing.
  - Account for confounding
  - Need appropriate corrections (e.g. False Discovery Rate)
  - Scalability to large cohorts
- Win: Large dataset allow to test modeling assumptions / fit better models
  - Inference of confounding structures
  - Not possible before large-scale hypothesis testing/large datasets
  - More power due to large datasets
  - Gain in power by joint analysis of multiple traits

N=10

ATGACCTG**A**AACTGGGGGGA**C**TGACGTG**G**AACGGT ATGACCTG**C**AACTGGGGGGA**C**TGACGTG**C**AACGGT ATGACCTG**C**AACTGGGGGGA**C**TGACGTG**C**AACGGT ATGACCTG**A**AACTGGGGGGA**T**TGACGTG**G**AACGGT ATGACCTG**C**AACTGGGGGGA**T**TGACGTG**C**AACGGT ATGACCTG**C**AACTGGGGGGA**T**TGACGTG**C**AACGGT



P=10<sup>6</sup>



### Population structure (genetic)





#### LINEAR MODEL









NOISE

 $oldsymbol{\psi} \sim \mathcal{N}\left( oldsymbol{0}, \sigma_{e}^{2} 
ight)$ 

### Flowering in A. thaliana





### Flowering in A. thaliana

### Multi-modal data integration across molecular layers



# Association genetics with high-dimensional molecular phenotypes



**EMBL-EBI** 

### ranslation

### challenge:

- multiple testing
- non-genetic sources of variation

# Association genetics with high-dimensional phenotypes



![](_page_12_Picture_2.jpeg)

## Linear mixed models to account for sample-to-sample covariance

![](_page_13_Figure_1.jpeg)

![](_page_13_Picture_3.jpeg)

### Confounding factors: genetic and non-genetic structure

![](_page_14_Figure_1.jpeg)

![](_page_14_Picture_3.jpeg)

## Confounding factors: genetic and non-genetic structure

Single marker genetic mapping

![](_page_15_Figure_2.jpeg)

Nature, Lappalainen et al. 2013

![](_page_15_Picture_4.jpeg)

### Confounding factors: genetic and non-genetic structure

>non-genetic (batch/env) >genetic confounding (population structure)

![](_page_16_Figure_2.jpeg)

![](_page_16_Picture_3.jpeg)

![](_page_16_Picture_4.jpeg)

# Extending linear mixed models: beyond single-SNP single phenotype analyses

- Statistical challenges in high-dimensional association genetics
  - Normalization and scaling of quantitative trail Fusi et al., Nat Comm (2014)
  - Accounting for epistasis and non-linear genetic interactions Stephan et al., Nat Comm (2015)
  - Joint modeling of multiple traits and variants

Casale et al., Nat Meth (2015)

![](_page_17_Picture_6.jpeg)

1

![](_page_18_Figure_1.jpeg)

![](_page_19_Figure_1.jpeg)

Listgarten et al. Bioinformatics (2013)

![](_page_20_Figure_1.jpeg)

![](_page_21_Figure_1.jpeg)

outcomes

Casale et al., Nat Meth (2015)

# mtSet: aggregation across traits and causal variants

![](_page_22_Figure_1.jpeg)

Challenge: Cubical scaling means such an algorithm is impractical for even moderately-size datasets!

### Efficient inference for large-scale GWAS

![](_page_23_Figure_1.jpeg)

(human chrom20, 3,975 set tests for 4 traits)

![](_page_23_Picture_4.jpeg)

![](_page_24_Figure_0.jpeg)

## Summary so far

- Linear mixed models enable accounting for sample heterogeneity in genetic analyses
  - genetic (population structure): reduces false positives
  - non-genetic (batch, environment): increases power
- Scalability to large datasets
- Joint modeling of multiple (correlated) traits and multiple causal variants

# Accounting for heterogeneity between individuals and single-cells

DNA

mRNA

proteins

organ-level phenotypes

transcription

translation

 $(\bar{y}_{1,\cdot})$ 

![](_page_26_Figure_1.jpeg)

### Statistical genetics

- ► Software (LIMIX)<sup>1</sup>
- Phenotype normalisation<sup>2</sup>
- Modelling epistatic relationships<sup>3</sup>
- Joint modelling of correlated traits<sup>4</sup>

### Molecular heterogeneity

Genetics of gene expression
 Causality<sup>5</sup>

ATGACCTGAAACTGGGGGACTGACGTGAACGG ATGACCTGCAACTGGGGGACTGACCTGCAACGGT ATGACCTGCAACTGGGGGACTGACGTGCAACGGT ATGACCTGAAACTGGGGGATGACGTGGCAACGG

- Human iPS biology
  Drosophila/yeast/plant genetics
  Cancer
- Drug susceptibility screens

![](_page_26_Picture_11.jpeg)

4.Casale & Rakitsch et al., Nat Meth (2015)5.Gagneur et al. PLoS Genet (2013)6.Buettner et al., Nat BioTech (2015)

![](_page_26_Picture_13.jpeg)

### single-cell heterogeneity

- Modelling heterogeneity in scRNA-Seq<sup>6</sup>
- Single cell DNA methylation profiling

![](_page_26_Picture_17.jpeg)

## Single-cell RNA-Seq

- Conventional RNA-Seq profiles are obtained from a pool of typically ~100,000+ cells.
- Using single-cell RNA-sequencing technologies, we can now assay RNA abundance in single cells.

- novel variation between cells: cell type composition, differentiation
- additional (confounding) expression heterogeneity: cell cycle, apoptosis, ...

![](_page_27_Picture_5.jpeg)

Fluidigm C1®

![](_page_27_Picture_7.jpeg)

# Cell cycle masks differentiation processes in single-cell RNA-Seq

![](_page_28_Figure_1.jpeg)

 Observed expression profiles do not enable recovering of the differentiation process.

EMBL-EBI

between cell cycle genes

and non-cycle genes

### Gene expression heterogeneity is not new...

![](_page_29_Figure_1.jpeg)

EMBL-EBI 🌡

## Single-cell latent variable model (scLVM)

- Random effect model for cell cycle effects. Two-stage approach:
  - 1. Estimate a cell-cell

![](_page_30_Picture_3.jpeg)

![](_page_30_Picture_4.jpeg)

Florian Kedar Buettner Natarajan

![](_page_30_Picture_6.jpeg)

Estimation of cell-cycle induced

![](_page_30_Picture_8.jpeg)

## Estimating the cell cycle covariance

+

cell cycle covariance

 $\delta_b \mathbf{I}$ 

- Reconstruct cell cycle from the observed expression data
- Use known annotated cell cycle gene set
- Employ latent variable modeling to reconstruct a cell cycle factor (X)

 $\mathbf{Y}_{\mathrm{cc}} \sim \left[ \begin{array}{c} \mathcal{N}(\mathbf{0} \mid \mathbf{0}) \right]$ 

g

![](_page_31_Picture_4.jpeg)

### **Technical noise requires special attention**

- Large proportions of technical variability due to low quantities of starting material
- Estimation of technical noise
  - Mean/variance fit from ERCC spike ins
  - Extrapolation to genome-wide genes
  - 7,073 highly variable genes

![](_page_32_Figure_6.jpeg)

Brennecke et al. 2013

![](_page_32_Picture_8.jpeg)

![](_page_32_Picture_9.jpeg)

## Decomposing sources of gene expression variation

- Variance decomposition of gene expression, considering
  - cell cycle (using estimated covariance)
  - residual biological variability
  - technical noise (estimated via spike-ins)

$$\mathbf{Y}_{g} = \boldsymbol{\mu} \mathbf{I} + \boldsymbol{\alpha} \mathbf{u}_{cc} + \boldsymbol{\delta}_{b} \mathbf{u}_{b} + \mathbf{u}_{n}$$

$$N(0, \boldsymbol{\mu}) N(0, \boldsymbol{\mu}) N(0, \boldsymbol{\mu})$$

$$N(0, \boldsymbol{\mu}) N(0, \boldsymbol{\mu}) N(0, \boldsymbol{\mu})$$

$$N(0, \boldsymbol{$$

### Model validation on mouse ESCs

 To test our model, we used single-cell RNA-Seq data generated from ~300 ES cells collected at different stages of the cell cycle

![](_page_34_Figure_2.jpeg)

- scLVM accurately estimates variability due to the cell cycle.
- Cell cycle effects are not visible on the model residuals.

## **Application to T-cell differentiation**

![](_page_35_Figure_1.jpeg)

- Focus on cells being differentiated in vitro from the naïve state towards the Th2 cell type
- 96 cells transcription profiled using the Fluidigm C1 system

![](_page_35_Picture_4.jpeg)

# Dissecting the sources of transcriptional variation

### Technical noise

For 27% of the genes, variation of expression can be entirely explained by the (technical) null variability.

### Cell-cycle

For 42% of the genes, >30% of the observed variance is explained by the cell cycle state.

![](_page_36_Figure_5.jpeg)

## The impact of cell cycle on gene-gene correlations

Gene-gene correlations (adjusted for cell cycle)

Gene-gene correlations (unadjusted)

GO.ID Term Annotated Significant Expected result1 G0:0006412 translation 416 55 6.49 8.0e-17 1 G0:0006414 translational elongation 45 13 0.70 1.2e-13 2 ribosomal small subunit assembly 10 GO:000028 6 0.16 2.8e-09 3 ADP biosynthetic process G0:0006172 8 5 0.12 4.8e-08 5 G0:0015986 17 6 ATP synthesis coupled proton transport 0.27 1.5e-07 > 500 G0:0006096 59 8 0.92 3.6e-06 6 glycolysis 92 G0:0006413 translational initiation 12 1.44 9.5e-06 21 GO:0001916 positive regulation of T cell mediated c... 5 0.33 1.5e-05 8 G0:0071353 cellular response to interleukin-4 22 5 0.34 1.9e-05 9 64Z 10 GO:0008284 positive regulation of cell proliferatio... Z8 10.02 Z.6e-05 11 GO:0000462 maturation of SSU-rRNA from tricistronic... 5 3 0.08 3.7e-05 12 GO:0015991 ATP hydrolysis coupled proton transport 25 5 0.39 3.7e-05 13 GO:0006662 glycerol ether metabolic process 13 4 0.20 3.7e-05 14 GO:0002474 antigen processing and presentation of p... 19 6 0.30 5.0e-05 15 G0:0042273 ribosomal large subunit biogenesis 14 4 0.22 5.2e-05

# Discrimination between differentiated and undifferentiated cells

![](_page_38_Figure_1.jpeg)

 After correction for cell heterogeneity, cells appear to separate better into two groups than without correction.

# Can we better tease apart the effect of cell cycle and differentiation ?

- scLVM also enables learning multiple latent factors
  - Genes annotated for cell cycle
  - Th2 differentiation marker genes

Extended variance component analysis

![](_page_39_Figure_5.jpeg)

![](_page_39_Picture_6.jpeg)

![](_page_39_Picture_7.jpeg)

# Can we better tease apart the effect of cell cycle and differentiation ?

### Th2 differentiation

928 genes with affected by the Th2 differentiation factor

- Th2/cell-cycle interaction
   200 genes with interaction effects
- Enriched for positive cell proliferation negative regulation of apoptosis

![](_page_40_Figure_5.jpeg)

![](_page_40_Picture_6.jpeg)

## The origin of transcriptome diversity?

![](_page_41_Picture_1.jpeg)

![](_page_41_Picture_2.jpeg)

Thierry

Voet

## Single-cell bisulfite sequencing (20 ESC cells)

![](_page_42_Figure_1.jpeg)

![](_page_42_Picture_2.jpeg)

### Parallel bs-seq & RNA-seq profiling in 21 serum ES cells

![](_page_43_Figure_1.jpeg)

![](_page_43_Picture_2.jpeg)

### Parallel bs-seq & RNA-seq profiling in 21 serum ES cells

![](_page_44_Figure_1.jpeg)

GRcm38 Chr12: 86361117 - 86521628 (161 kbp)

![](_page_44_Picture_3.jpeg)

### Parallel bs-seq & RNA-seq profiling in 21 serum ES cells

![](_page_45_Figure_1.jpeg)

![](_page_45_Picture_2.jpeg)

### Conclusions

- Latent variable models can effectively account for gene expression heterogeneity & confounding
- (e)QTL analysis
  - population structure & env. /technical confounding to improve power
- Single-cell RNA-seq analysis
  - a small number of genes with known cell cycle annotation is sufficient to estimate a cell covariance due to cell cycle
  - more compact gene-gene correlations
  - detection of genes with interactions involving multiple biological processes
  - Parallel bs-seq/RNA-seq profiling in the same cells reveals associations between methylation and transcriptome variation

![](_page_46_Picture_9.jpeg)

### Acknowledgments

<u>group</u> Barbara Rakitsch Florian Buettner Christof Angermüller Paolo Casale Helena Kilpinen Amelie Baud Danilo Horta

Johannes Stephan Bogdan Mirauta Kate Howell Fatemeh Ghavidel <u>EBI/Sanger</u> John Marioni Sarah Teichmann Thierry Voet

### Kedar Natarajan

Valentina Proserpio Antonio Scialdone Iain Macaulay

Babraham Inst. Wolf Reik Gavin Kelsey Heather Lee Stephen Clark <u>Helmholtz Munich</u> Fabian Theis

Microsoft Research Nicolo Fusi Christoph Lippert

University of Sheffield Neil Lawrence

> Postdoc opportunities: Single-cell genomics Statistical genetics

Mixed model software: <u>https://github.com/PMBio/limix</u>

scLVM: https://github.com/PMBio/scLVM

![](_page_47_Picture_13.jpeg)

![](_page_47_Picture_14.jpeg)

![](_page_47_Picture_15.jpeg)

![](_page_47_Picture_16.jpeg)