

eQTL analysis – an approach with Bioconductor

Vincent Carey

Channing Division of Network Medicine

Harvard Medical School

Road map

- Background on scope and implications of eQTL studies
- Data structures suitable for current studies (RNA-seq and genome-wide DNA sequencing/imputation [SNP only])
- Demonstration with Bioconductor
 - RNA-seq FPKM in geuvPack (special for this course)
 - 1000 Genomes VCF in an Amazon S3 bucket
 - “Batch effect” correction and analysis

Genetic analysis of genome-wide variation in human gene expression

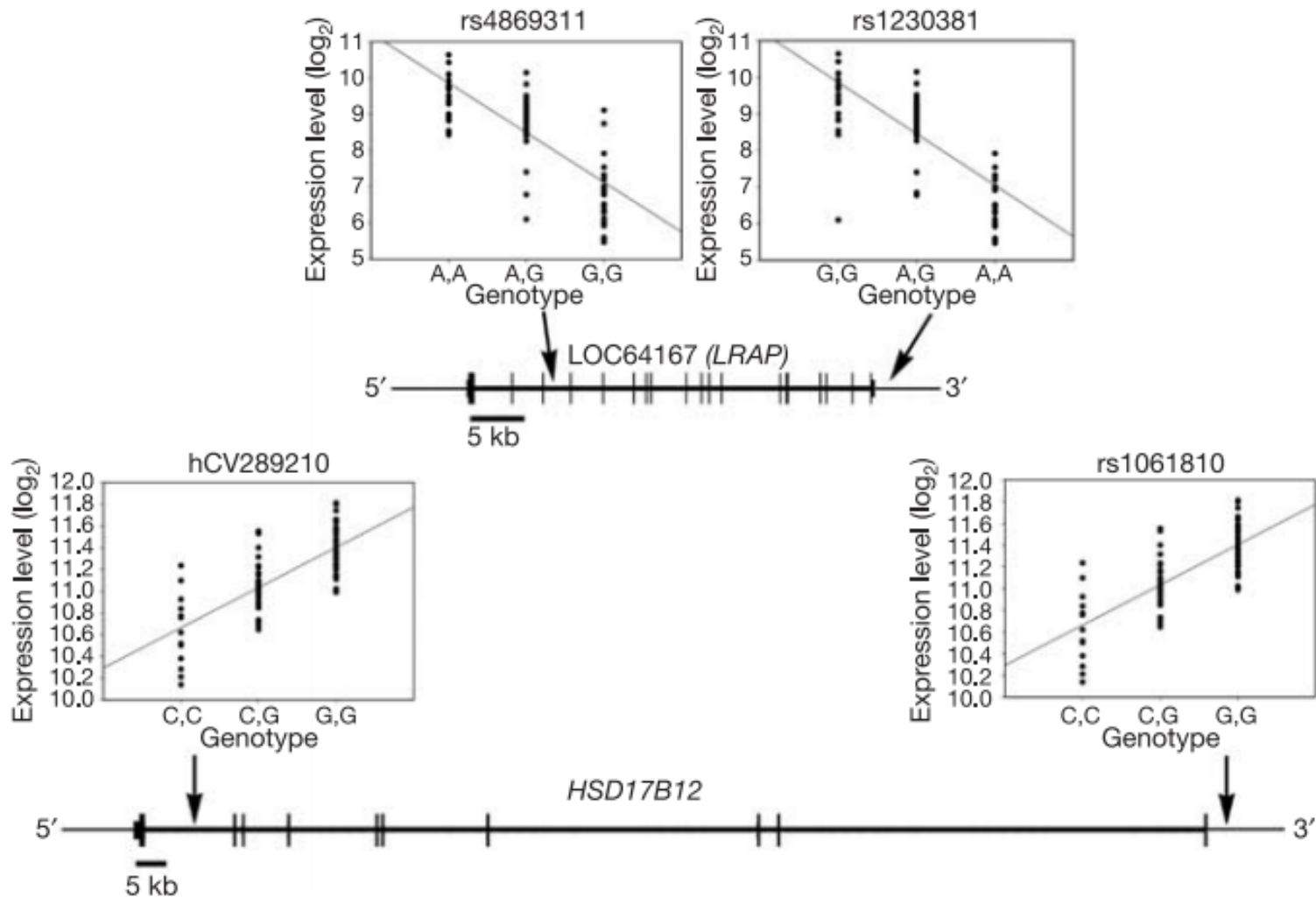
Michael Morley^{1,3*}, Cliona M. Molony^{2*}, Teresa M. Weber^{1,3}, James L. Devlin², Kathryn G. Ewens², Richard S. Spielman² & Vivian G. Cheung^{1,2,3}

¹Department of Pediatrics and ²Department of Genetics, University of Pennsylvania,

³The Children's Hospital of Philadelphia, Philadelphia, Pennsylvania 19104, USA

*These authors contributed equally to this work

Natural variation in gene expression is extensive in humans and other organisms, and variation in the baseline expression level of many genes has a heritable component. To localize the genetic determinants of these quantitative traits (expression phenotypes) in humans, we used microarrays to measure gene expression levels and performed genome-wide linkage analysis for expression levels of 3,554 genes in 14 large families. For approximately 1,000 expression phenotypes, there was significant evidence of linkage to specific chromosomal regions. Both *cis*- and *trans*-acting loci regulate variation in the expression levels of genes, although most act *in trans*. Many gene expression phenotypes are influenced by several genetic determinants. Furthermore, we found hotspots of transcriptional regulation where significant evidence of linkage for several expression phenotypes (up to 31) coincides, and expression levels of many genes that share the same regulatory region are significantly correlated. The combination of microarray techniques for phenotyping and linkage analysis for quantitative traits allows the genetic mapping of determinants that contribute to variation in human gene expression.

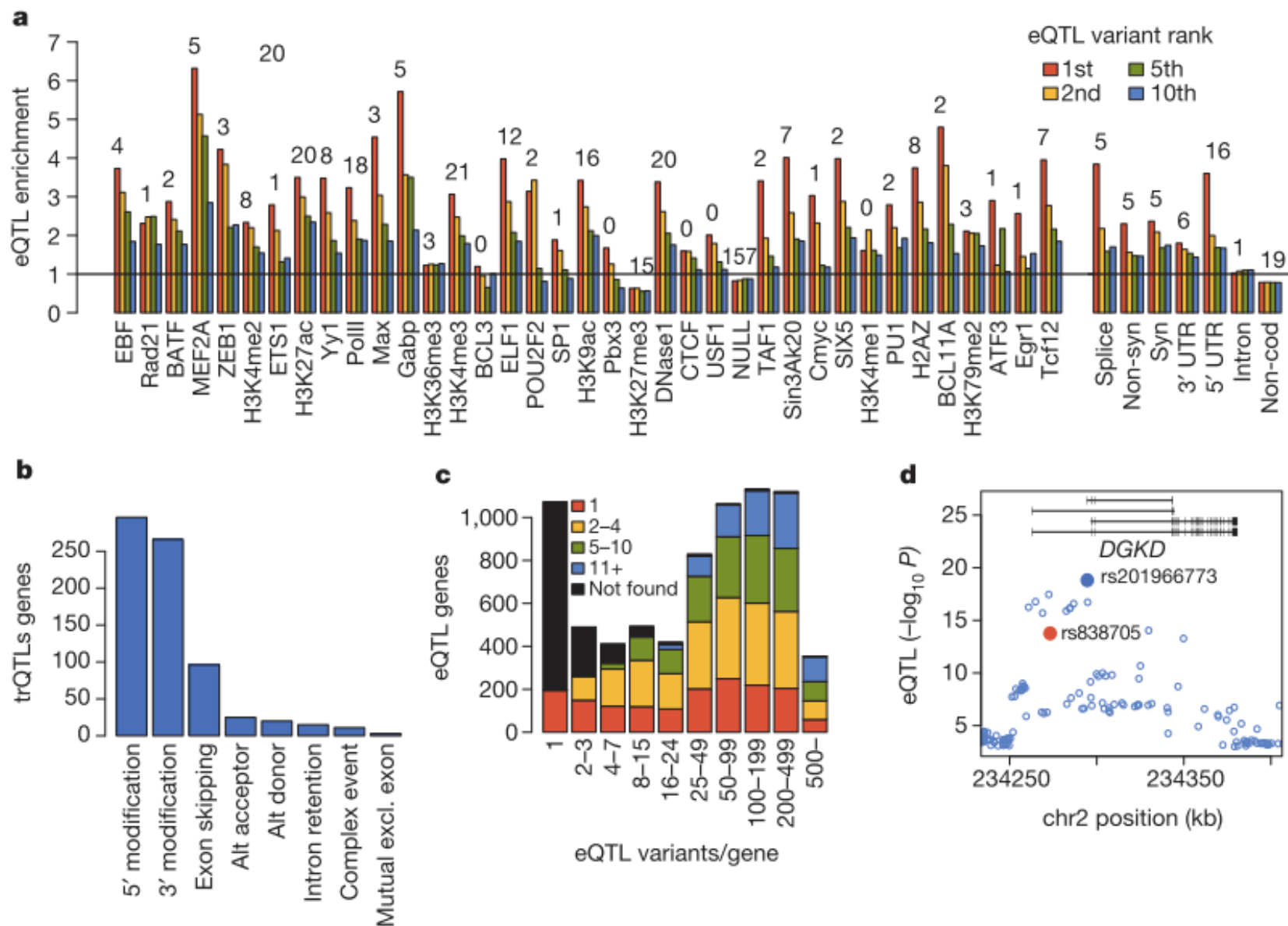


Association of expression phenotype of LOC64167 and *HSD17B12* on nearby SNPs. For LOC64167, the distance between SNP markers rs4869311 and rs1230381 is 172 kb. The distance between hCV289210 and rs1061810 is 172 kb.

Transcriptome and genome sequencing uncovers functional variation in humans

Tuuli Lappalainen^{1,2,3}, Michael Sammeth^{4,5,6,7,†*}, Marc R. Friedländer^{5,6,7,8*}, Peter A. C. 't Hoen^{9*}, Jean Monlong^{5,6,7*}, Manuel A. Rivas^{10*}, Mar González-Porta¹¹, Natalja Kurbatova¹¹, Thasso Griebel⁴, Pedro G. Ferreira^{5,6,7}, Matthias Barann¹², Thomas Wieland¹³, Liliana Greger¹¹, Maarten van Iterson⁹, Jonas Almlöf¹⁴, Paolo Ribeca⁴, Irina Pulyakhina⁹, Daniela Esser¹², Thomas Giger¹, Andrew Tikhonov¹¹, Marc Sultan¹⁵, Gabrielle Bertier^{5,6}, Daniel G. MacArthur^{16,17}, Monkol Lek^{16,17}, Esther Lizano^{5,6,7,8}, Henk P. J. Buermans^{9,18}, Ismael Padioleau^{1,2,3}, Thomas Schwarzmayr¹³, Olof Karlberg¹⁴, Halit Ongen^{1,2,3}, Helena Kilpinen^{1,2,3}, Sergi Beltran⁴, Marta Gut⁴, Katja Kahlem⁴, Vyacheslav Amstislavskiy¹⁵, Oliver Stegle¹¹, Matti Pirinen¹⁰, Stephen B. Montgomery^{1,†}, Peter Donnelly¹⁰, Mark I. McCarthy^{10,19}, Paul Flicek¹¹, Tim M. Strom^{13,20}, The Geuvadis Consortium[‡], Hans Lehrach^{15,21}, Stefan Schreiber¹², Ralf Sudbrak^{15,21,†}, Ángel Carracedo²², Stylianos E. Antonarakis^{1,2}, Robert Häsler¹², Ann-Christine Syvänen¹⁴, Gert-Jan van Ommen⁹, Alvis Brazma¹¹, Thomas Meitinger^{13,20,23}, Philip Rosenstiel¹², Roderic Guigó^{5,6,7}, Ivo G. Gut⁴, Xavier Estivill^{5,6,7,8} & Emmanouil T. Dermitzakis^{1,2,3}

Genome sequencing projects are discovering millions of genetic variants in humans, and interpretation of their functional effects is essential for understanding the genetic basis of variation in human traits. Here we report sequencing and deep analysis of messenger RNA and microRNA from lymphoblastoid cell lines of 462 individuals from the 1000 Genomes Project—the first uniformly processed high-throughput RNA-sequencing data from multiple human populations with high-quality genome sequences. We discover extremely widespread genetic variation affecting the regulation of most genes, with transcript structure and expression level variation being equally common but genetically largely independent. Our characterization of causal regulatory variation sheds light on the cellular mechanisms of regulatory and loss-of-function variation, and allows us to infer putative causal variants for dozens of disease-associated loci. Altogether, this study provides a deep understanding of the cellular mechanisms of transcriptome variation and of the landscape of functional variants in the human genome.



transcriptome QTLs. **a**, Enrichment of EUR exon eQTLs in 100 permutations for the first, second, fifth and tenth best associating variant per gene, relative to a matched null set of variants denoted by the horizontal line. The numbers are $-\log_{10}(P)$ values of a Fisher test between the

caused by transcript ratio QTLs. **c**, The rank of the best Omni 2.5 variant for the significant EUR eQTL variants per gene. **d**, The *DGKD* gene located on chromosome 2. An intronic SNP rs838705 is associated with calcium levels (red dot). The eQTL variant 21 kb downstream (blue dot) is a very likely causal variant

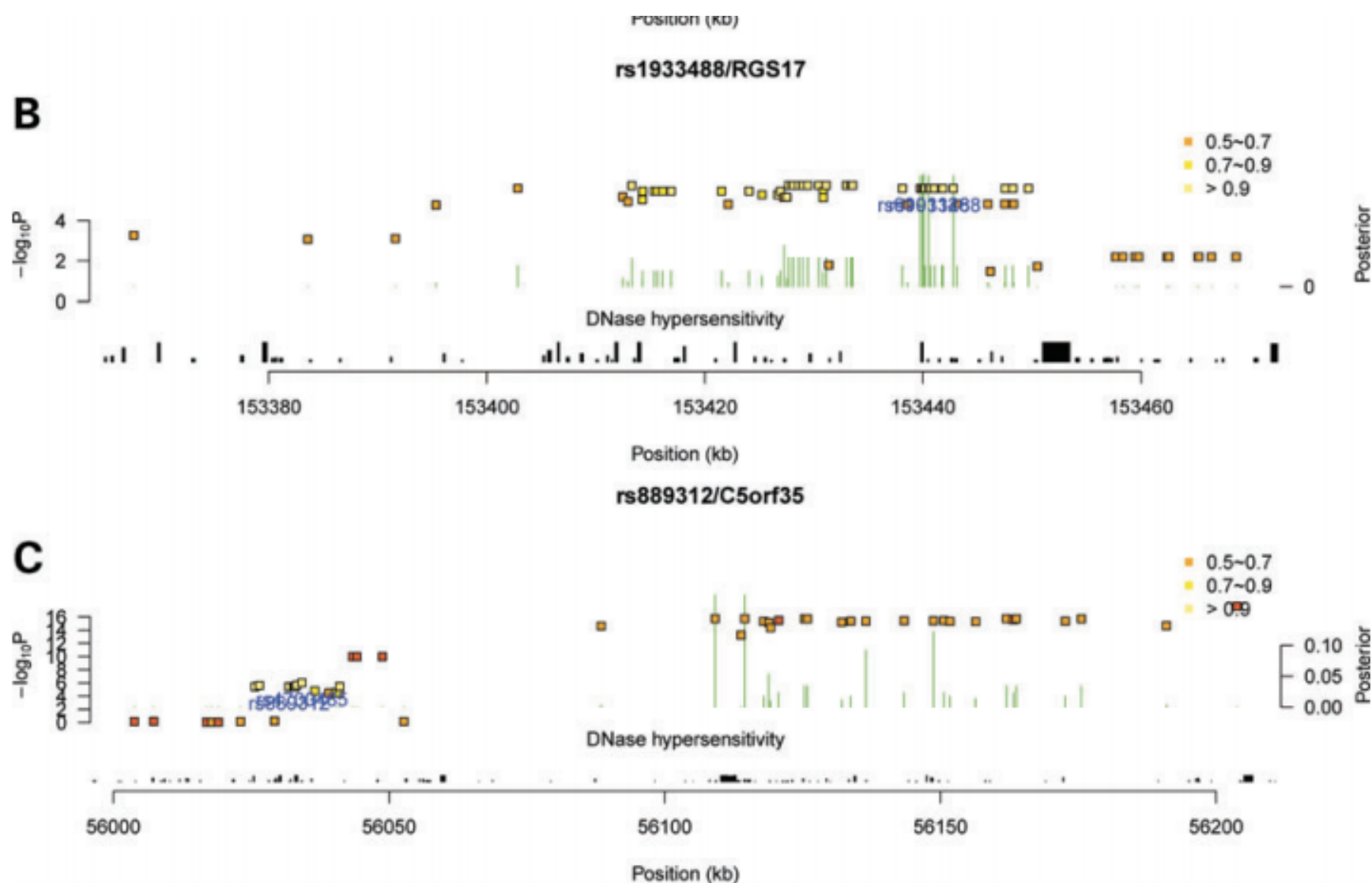
Expression QTL-based analyses reveal candidate causal genes and loci across five tumor types

Qiyuan Li^{1,2,3}, Alexander Stram⁴, Constance Chen⁵, Siddhartha Kar⁶, Simon Gayther⁷, Paul Pharoah⁶, Christopher Haiman⁴, Barbara Stranger⁸, Peter Kraft⁵ and Matthew L. Freedman^{1,3,*}

¹Department of Medical Oncology, The Center for Functional Cancer Epigenetics, Dana Farber Cancer Institute, Boston, MA, USA, ²Medical College of Xiamen University, Xiamen, China, ³Program in Medical and Population Genetics, The Broad Institute, Cambridge, MA, USA, ⁴University of Southern California, Los Angeles, CA, USA, ⁵Department of Epidemiology, Harvard School of Public Health, Boston, MA, USA, ⁶Strangeways Research Laboratory, University of Cambridge, Cambridge, UK, ⁷Department of Preventive Medicine, Keck School of Medicine, University of Southern California, Log Angeles, CA, USA and ⁸Institute for Genomics and Systems Biology, University of Chicago, Chicago, IL, USA

Received February 17, 2014; Revised April 17, 2014; Accepted May 6, 2014

The majority of trait-associated loci discovered through genome-wide association studies are located outside of known protein coding regions. Consequently, it is difficult to ascertain the mechanism underlying these variants and to pinpoint the causal alleles. Expression quantitative trait loci (eQTLs) provide an organizing principle to address both of these issues. eQTLs are genetic loci that correlate with RNA transcript levels. Large-scale data sets such as the Cancer Genome Atlas (TCGA) provide an ideal opportunity to systematically evaluate eQTLs as they have generated multiple data types on hundreds of samples. We evaluated the determinants of gene expression (germline variants and somatic copy number and methylation) and performed *cis*-eQTL ana-



ation of the fine-mapping candidates of three cancer risk loci based on an integrated posterior probability combining association presents a correlated germline variant of the initially reported risk locus (labeled by blue text); the height of the points corresponds to $-\log_{10} P$ values; the DNaseI HS scores are shown beneath the posterior; the green bars show the posterior probabilities. (A) Xp11.232/rs1044396 locus for prostate cancer; (B) 6q25.2/rs1933488 with *RGS17* in prostate cancer and (C) 5q11.2/rs889312 with *C5orf35* in ER-positive breast cancer.

Common Genetic Variants Modulate Pathogen-Sensing Responses in Human Dendritic Cells

Mark N. Lee,^{1,2,3*} Chun Ye,^{1*} Alexandra-Chloé Villani,^{1,2} Towfique Raj,^{1,2,4} Weibo Li,^{1,3} Thomas M. Eisenhaure,^{1,3} Selina H. Imboywa,² Portia I. Chipendo,² F. Ann Ran,^{1,5,6,7,8} Kamil Slowikowski,⁹ Lucas D. Ward,^{1,10} Khadir Raddassi,¹¹ Cristin McCabe,^{1,4} Michelle H. Lee,² Irene Y. Frohlich,² David A. Hafler,⁸ Manolis Kellis,^{1,10} Soumya Raychaudhuri,^{1,2,12,13} Feng Zhang,^{6,7,8} Barbara E. Stranger,^{14,15} Christophe O. Benoist,² Philip L. De Jager,^{1,2,4} Aviv Regev,^{1,16,17} † Nir Hacohen^{1,2,3} †

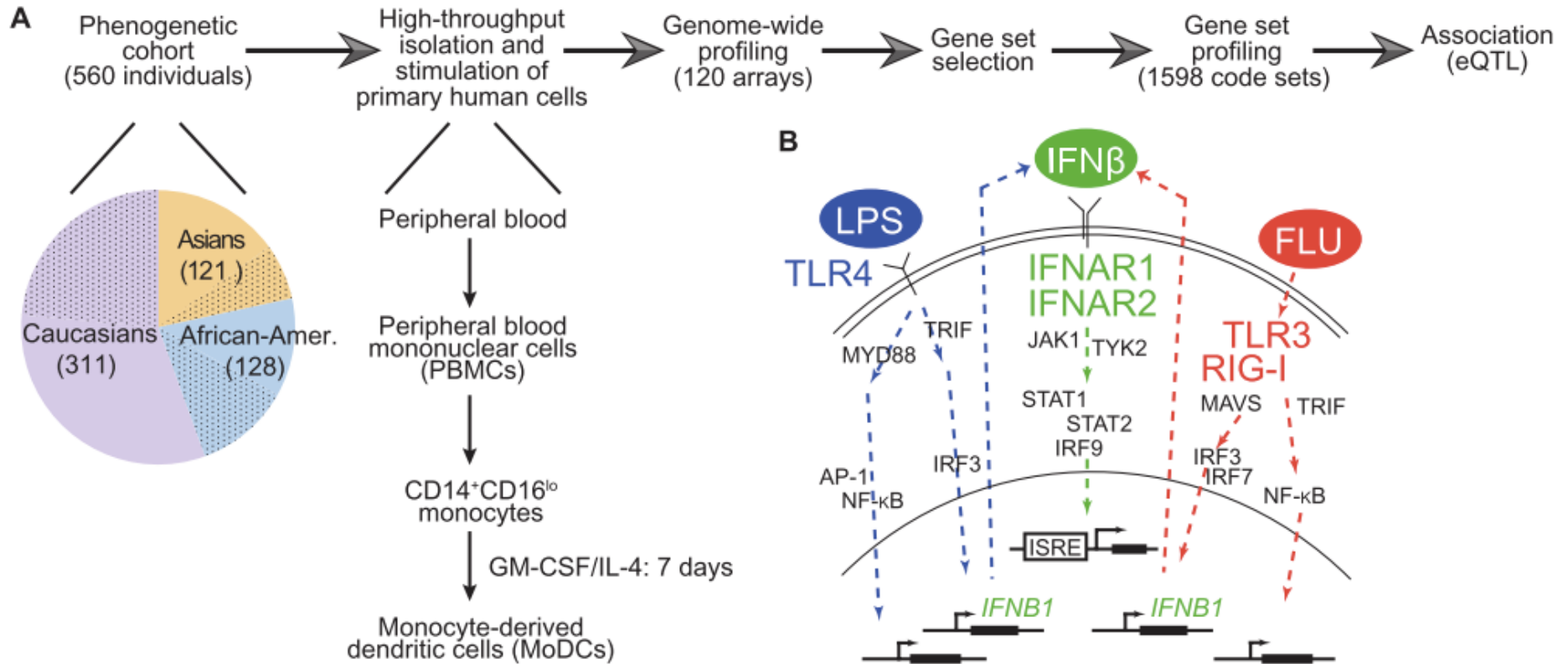
Little is known about how human genetic variation affects the responses to environmental stimuli in the context of complex diseases. Experimental and computational approaches were applied to determine the effects of genetic variation on the induction of pathogen-responsive genes in human dendritic cells. We identified 121 common genetic variants associated in cis with variation in expression responses to *Escherichia coli* lipopolysaccharide, influenza, or interferon- β (IFN- β). We localized and validated causal variants to binding sites of pathogen-activated STAT (signal transducer and activator of transcription) and IRF (IFN-regulatory factor) transcription factors. We

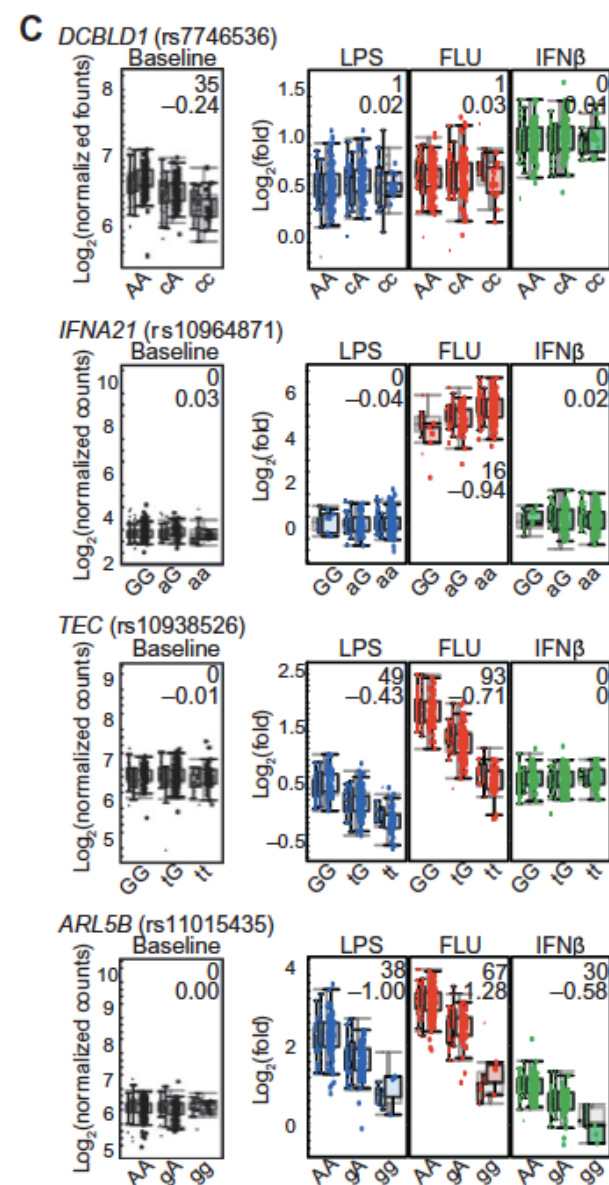
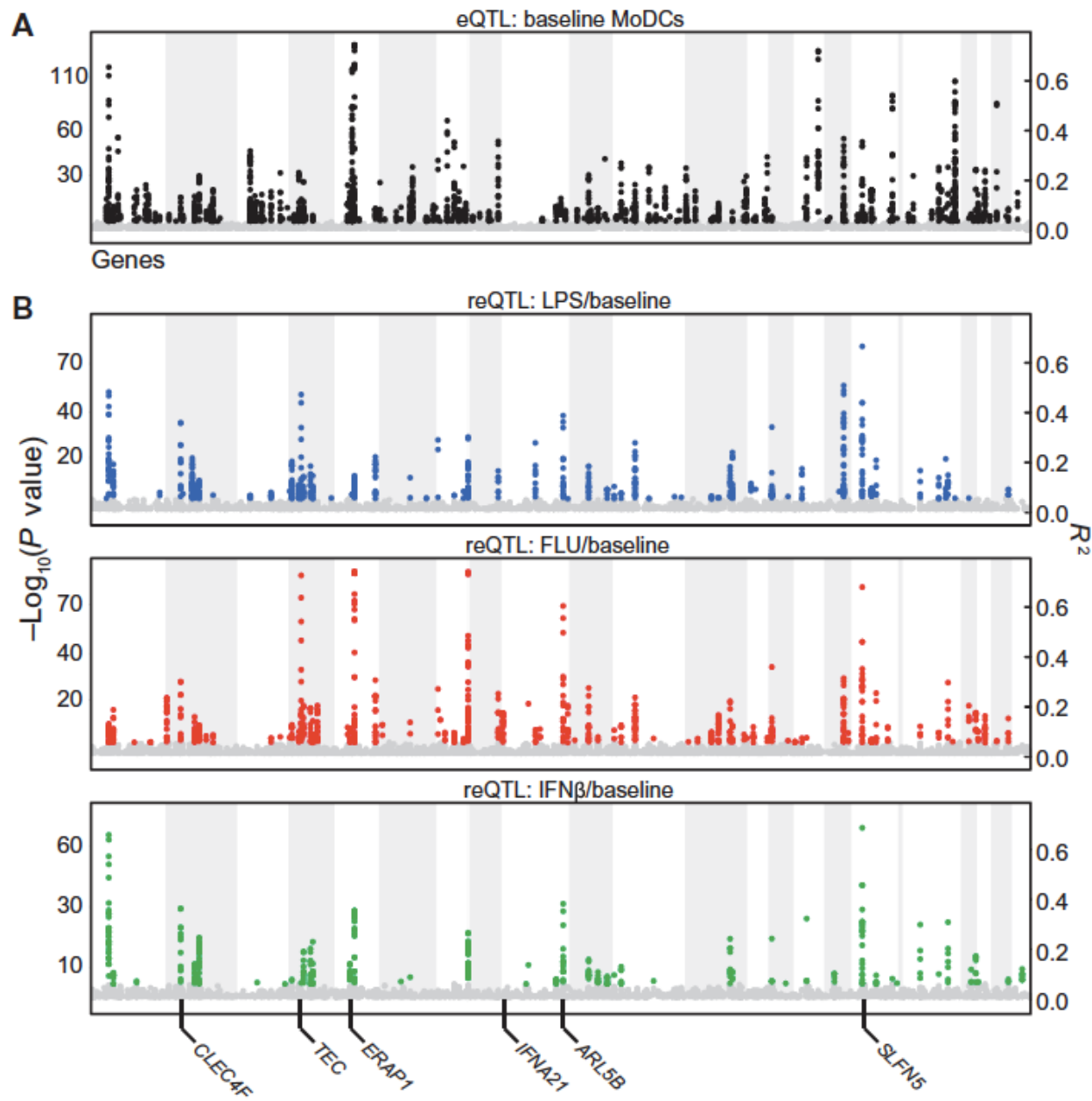
cleus to ind
including ir
gages the ty
the expressi
Genetic stu
iants near r
risk of diffc
DCs also p
immune res
eases (11–1
wide associa
(14–17), esp

Results

Assessing t Pathogen S

We develop
computatio
variability i
this variabi
First, we o
to isolate pri
human bloo





Summary

- 2004: Microarray methods can expose genetic variants associated with expression variation
- 2011-2013: RNA-seq applications expose genetic sources of detailed variations in transcription (e.g., alternate splicing, allele-specific expression)
- 2014+: Tissue- and environment-specificity of eQTL associations to the fore
- Upshot – more eQTL experiments, reference data, detailed annotation for integrative use

Major concerns of this tutorial

- How can I do my own eQTL inferences?
- How can I use eQTL inferences obtained elsewhere?
- These questions imply concern with
 - Efficient representation of inputs to eQTL analyses
 - Efficient computation of eQTL inferences
 - Efficient representation of eQTL inferences

Difficult feature of eQTL analyses

- For the comprehensive analysis of trans-eQTL associations, we are considering $|G| \times |S|$ tests, which ranges into the hundreds of billions
 - You are *creating* big data even if what you are starting with is modest in volume
 - *After you've done the tests*, interactive general query resolution into tables with millions of records requires computing methods that are somewhat unusual relative to classical statistical computing
 - You will want to use scalable computing strategies as the standard tests are embarrassingly parallel
 - A typical approach (see MatrixEQTL) is to discard information on apparently insignificant associations, given a prior threshold, as soon as possible in the analysis

Representing the input data on Expression

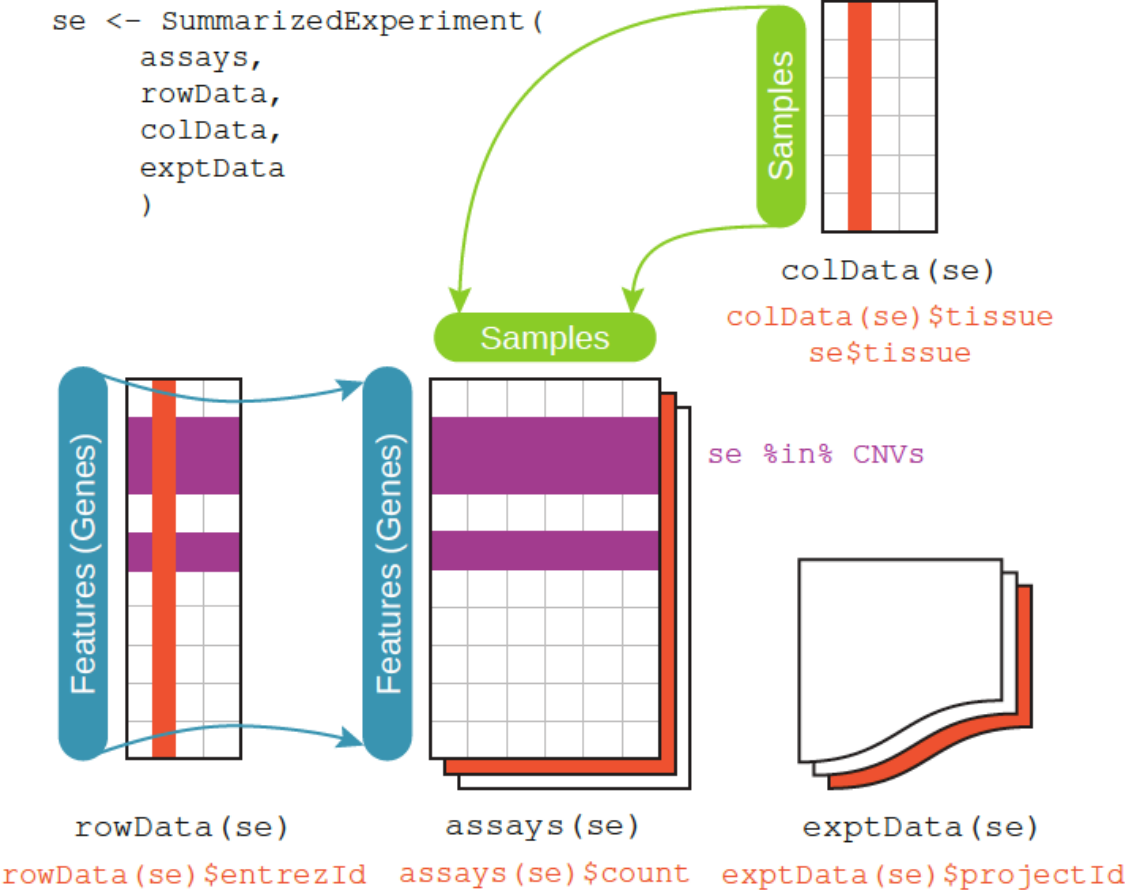



Figure 4: Elements from which *SummarizedExperiment* objects are composed.

The “FPKM” representation of expression in GEUVADIS samples

```
Console ~/   
R » library(geuvPack)  
R » data(geuFPKM)  
R » geuFPKM  
class: SummarizedExperiment  
dim: 23722 462  
exptData(2): MIAME constrHist  
assays(1): exprs  
rownames(23722): ENSG00000152931.6 ENSG00000183696.9 ...  
      ENSG00000257337.1 ENSG00000177494.5  
rowData metadata column names(18): source type ... tag ccdsid  
colnames(462): HG00096 HG00097 ... NA20826 NA20828  
colData names(0):  
R » exptData(geuFPKM)$MIAME  
Experiment data  
  Experimenter name: Lappalainen T  
  Laboratory: NA  
  Contact information:  
  Title: Transcriptome and genome sequencing uncovers functional variation in humans.  
  URL:  
  PMIDs: 24037378
```

Moderately self-documenting; currently lacks ‘sample info’

Quiz questions

- Q1. Use `dim(assay(geuFPKM))` to determine G (number of 'genes') and N (number of samples).
- Q2. Use `rowData(geuFPKM)`... to determine the name of the gene measured in the first row of the assay data in `geuFPKM`. Note the type of gene.
- Q3. Tabulate the gene types assayed in the experiment. Hint – use `mcols()`
- Q4. Obtain a histogram of the distribution of expression values for this gene. Is a Gaussian model reasonable for this sample?

Q5. How can you verify the coordinates of DGKD using information in geuFPKM?

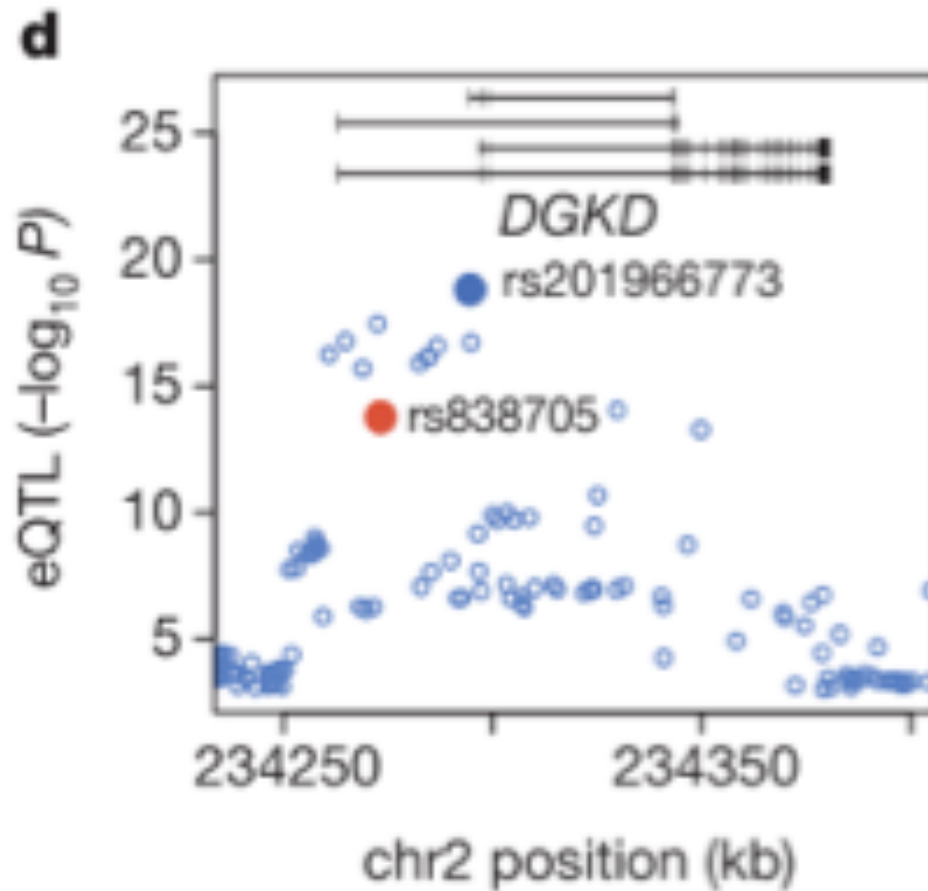


Figure 2d from the Geuvadis paper, PMID 24037378

Tools for approximate reproduction of Figure 2d

- We find the GEUVADIS quantification of DGKD in geuFPKM
- http://1000genomes.s3.amazonaws.com/release/20110521/ALL.chr2.phase1_release_v3.20101123.snps_indels_svs.genotypes.vcf.gz is the URL of a Tabix-indexed VCF with genotype calls from 1000 genomes
- GGtools cisAssoc() will use these two resources to compute association tests

Conceptual setup

- Perform an adjustment to expression data to remove technical variation
 - GEUVADIS used “PEER” to remove 10 latent factors; we’ll use 10 PC ... should be followed by sensitivity analysis; SVA is also of interest
- Compute association statistics of interest (in this case, we’ll use a cis radius of 50k)
- Correct p-values for multiple testing (we’ll use parametric testing; see the eQTL workflow for nonparametric approach)

Obtain cis eQTL for DGKD and one other gene on chr2

```
R » load("bfac.rda")          bfac are 10 PC from expression FPKM
R » data(geuFPKM)
R » colData(geuFPKM) = cbind(colData(geuFPKM), DataFrame(bfac))
R » source("buckpath.R")
R » cd = cisAssoc(geuFPKM[c(14894,42)], TabixFile(buckpath(2)), lbmaf=.05,
  cisradius=50000, rhs=~PC1+PC2+PC3+PC4+PC5+PC6+PC7+PC8+PC9+PC10)
```

Output of cisAssoc (first 3 records)

```
R » options(digits=3)
```

```
R » cd[1:3]
```

```
GRanges with 3 ranges and 10 metadata columns:
```

	seqnames	ranges	strand	paramRangeID	REF	ALT
	<Rle>	<IRanges>	<Rle>	<factor>	<DNAStringSet>	<CharacterList>
[1]	2	[234214947, 234214947]	*	ENSG00000077044.5	T	C
[2]	2	[234215226, 234215226]	*	ENSG00000077044.5	G	T
[3]	2	[234217249, 234217249]	*	ENSG00000077044.5	T	A

	chisq	permScore_1	permScore_2	permScore_3	snp	MAF	probeid
	<numeric>	<numeric>	<numeric>	<numeric>	<character>	<numeric>	<factor>
[1]	0.0046	0.4842	0.72	2.03	rs7588509	0.1675	ENSG00000077044.5
[2]	0.1487	0.0326	2.55	1.80	rs6719241	0.0618	ENSG00000077044.5
[3]	0.4409	0.0102	2.90	1.05	rs78571834	0.1425	ENSG00000077044.5

```
---
```

```
seqlengths:
```

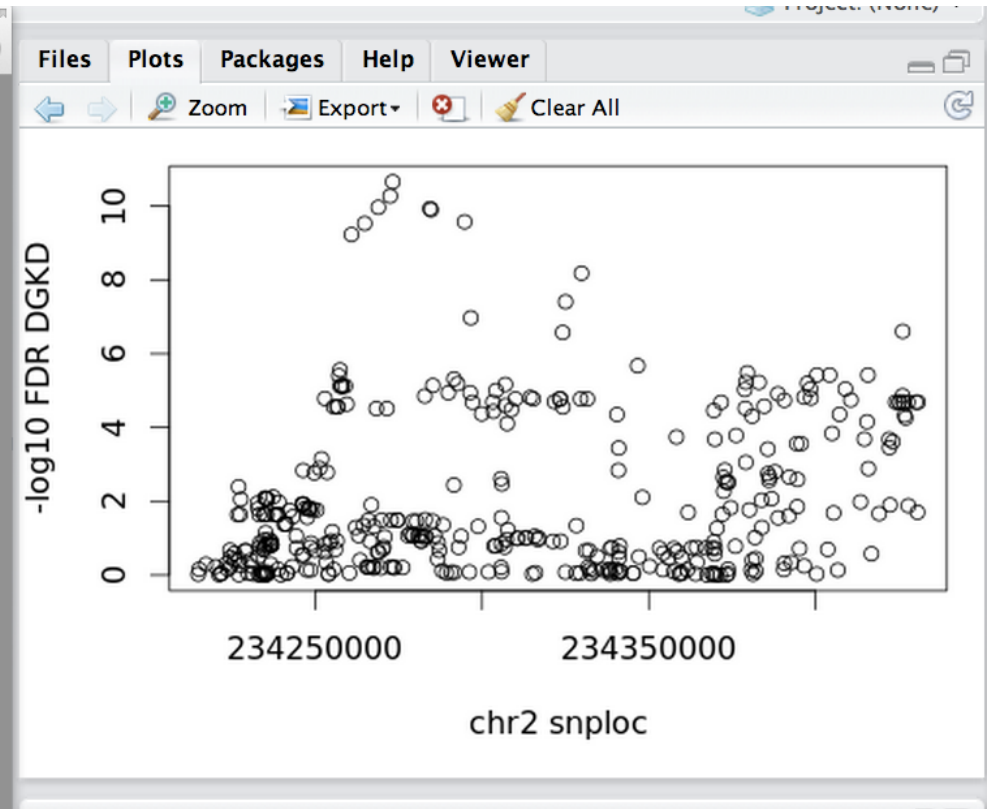
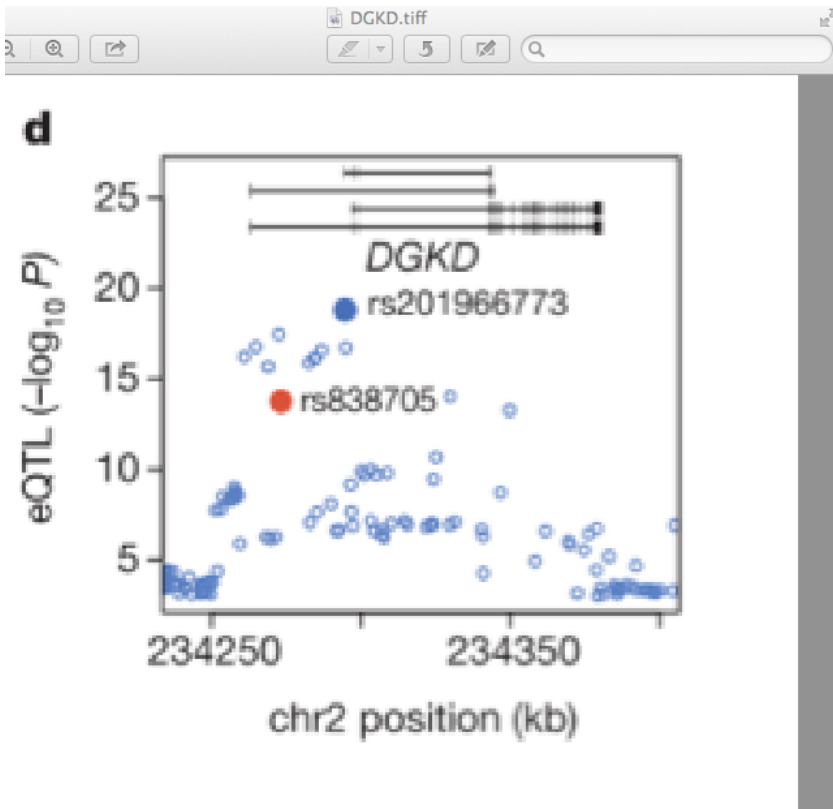
```
2
```

```
NA
```

Creating the manhattan plot

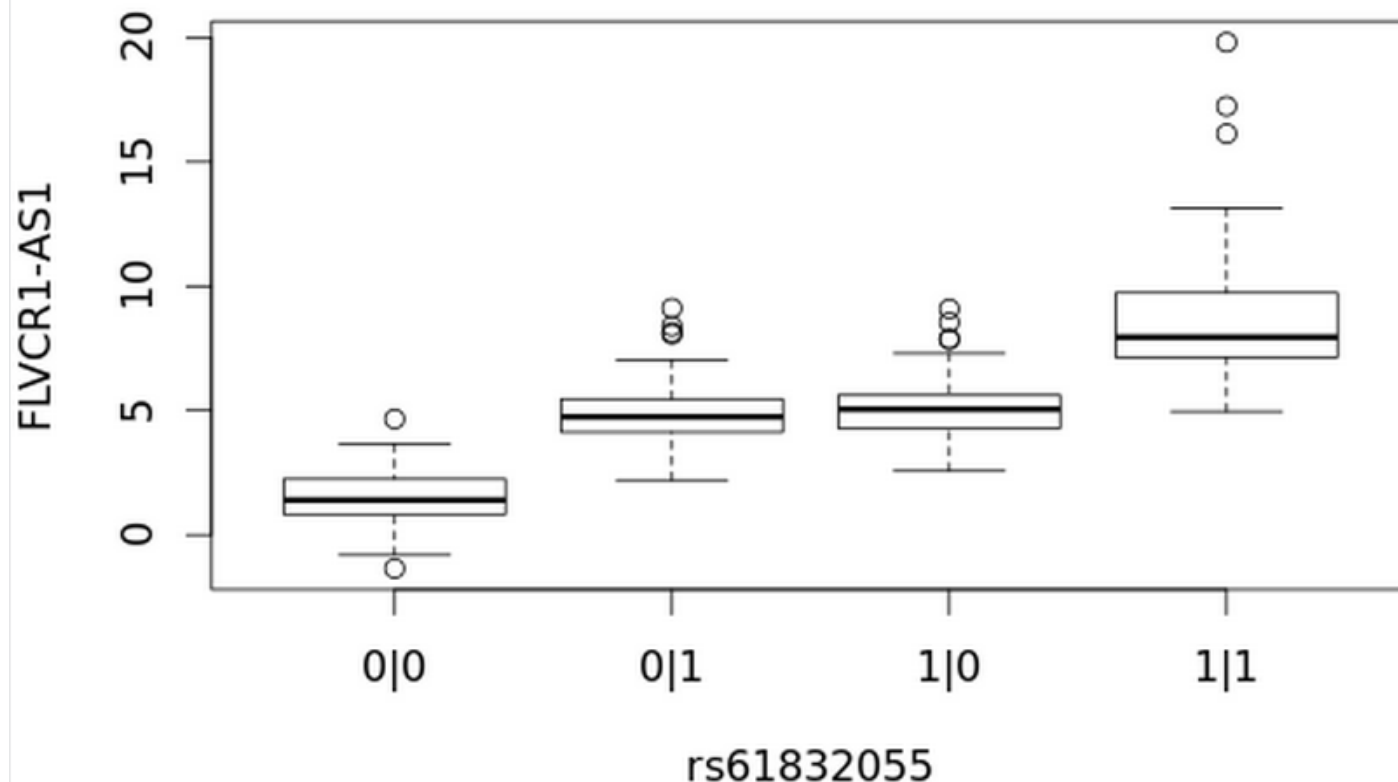
```
R » pv = 1-pchisq(cd$chisq,1)
R » library(multtest)
R » adp = mt.rawp2adjp(pv, proc="BY")
R » adpo = adp$adjp[order(adp$index)]
R » plot(start(cd)[1:404], -log10(adpo)[1:404])
R » plot(start(cd)[1:404], -log10(adpo)[1:404],
+   xlab="chr2 snploc", ylab="-log10 FDR DGKD"
+ )
```

Roughly in agreement...



Broader question: SNP associated with abundance of lincRNA transcripts

- Q1: create the subset of geuFPKM confined to genes assigned to type 'lincRNA'
- Q2: confine the subset to genes on chr1
- Q3: create the TabixFile reference to the 1000 genomes VCF for chr1
- Q4: use cisAssoc and plotOne to create the visualization



Creating the test statistics for FLVCR1- ASV

```
data(geuFPKM)
type = mcols(geuFPKM)$gene_type
islnc = which(type=="lincRNA")
length(islnc)
gl1 = geuFPKM[islnc,]
gl1 = gl1[which(seqnames(gl1)== "chr1"),]
t1 = TabixFile(buckpath(1))
seqlevelsStyle(gl1) = "NCBI"
library(VariantAnnotation)
library(BiocParallel)
cl1 = cisAssoc(gl1[3:4,], t1, lbmaf=.05, cisradius=50000)
```

To create the plot

- Determine the range associated with the SNP of interest
- Check the args of plotOne and supply the elements you've used to do the testing with the SNP range

Summary

- Genome-wide searches for cis-eQTL are easily carried out with
 - RNA-seq data in a SummarizedExperiment
 - Genotype data in a tabix-indexed VCF
 - Ggtools cisAssoc, with various approaches to FDR computation available
- See the ‘eQTL workflow’ at Bioconductor for more details on nonparametric inference, sensitivity analysis, and functional assessment