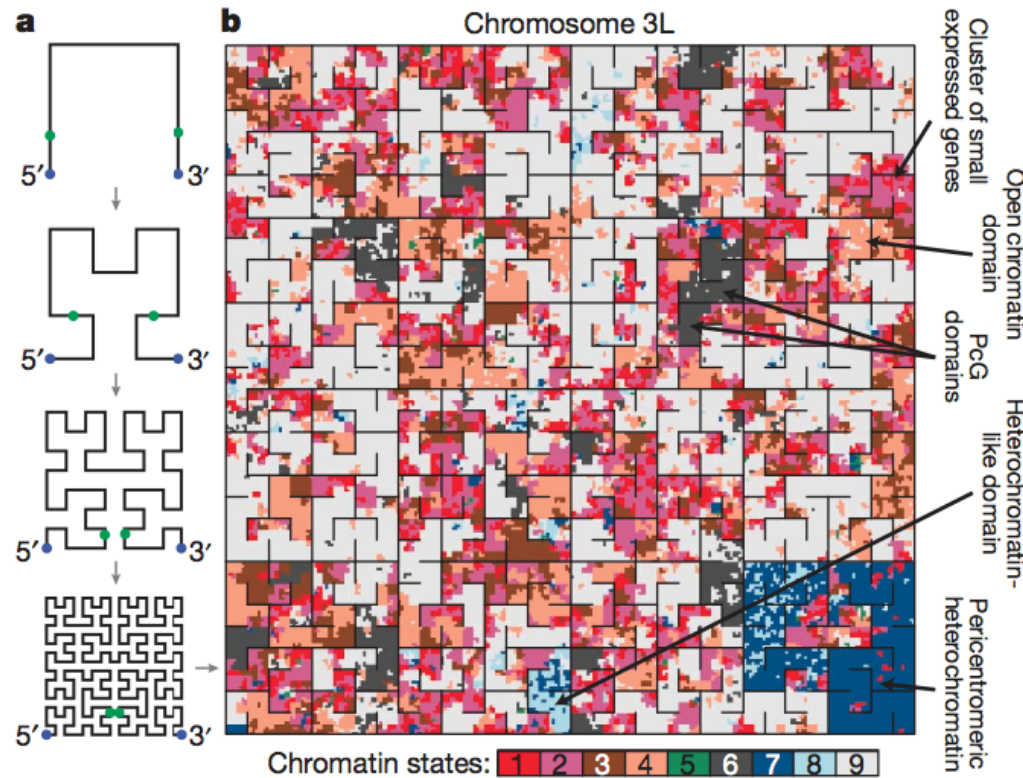


# Data structures and methods for some integrative analyses

Vince Carey PhD

Harvard Medical School

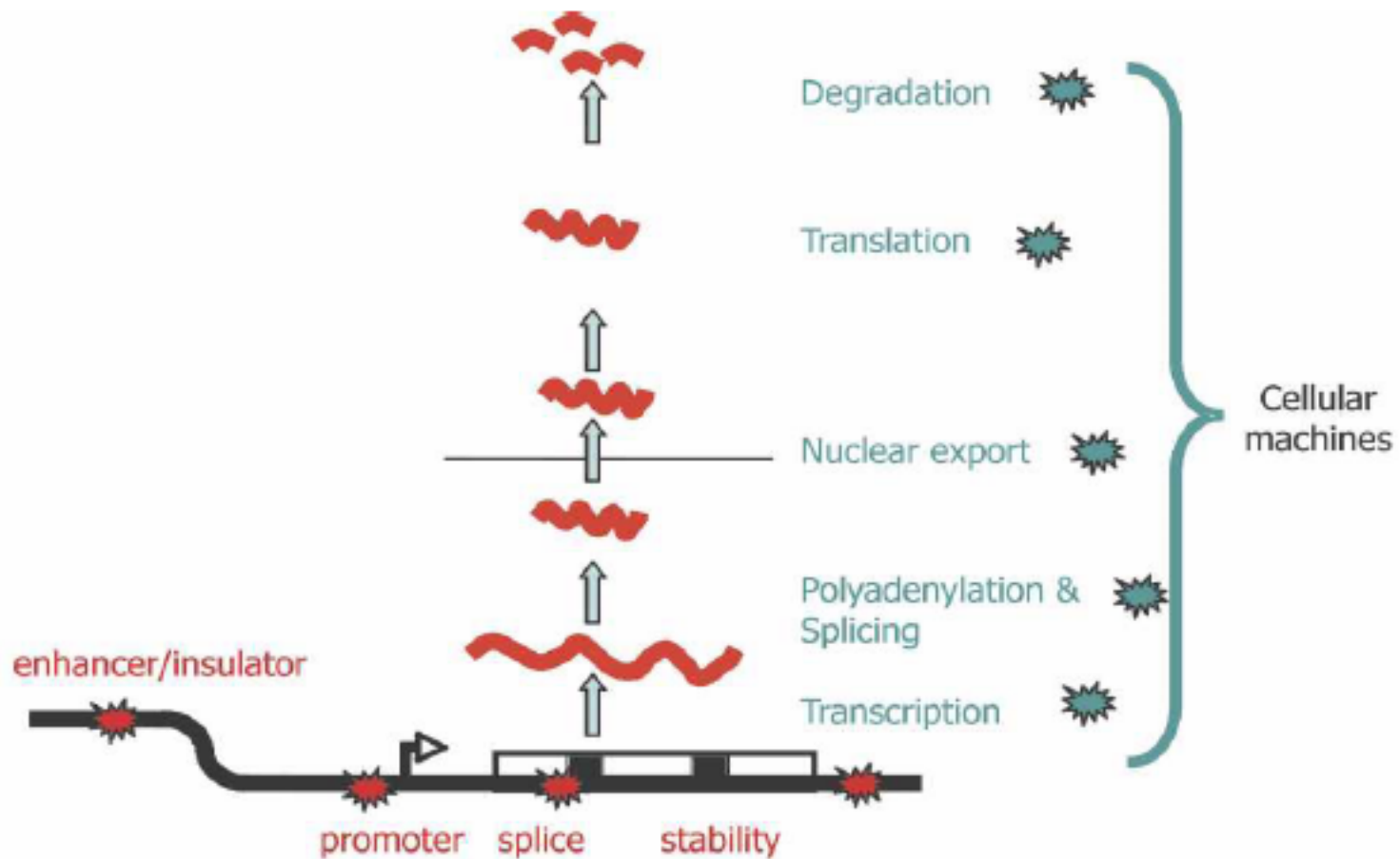
# The discreet charm of CSAMA listeners



So far four audience members have tapped me on the shoulder to help me understand fly chromosome structure ... glad to have your support!

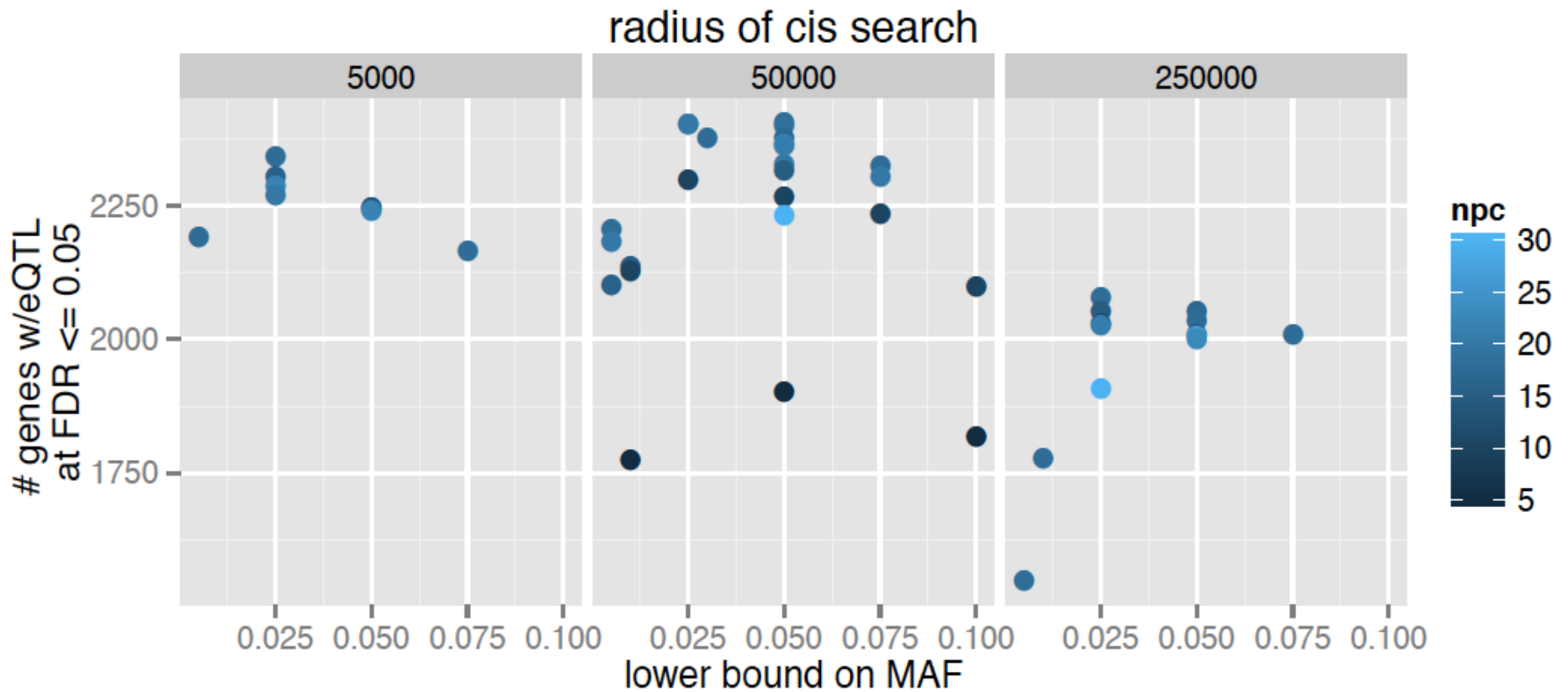
## 3 x 15'

- eQTL: sensitivity analysis, removal of extraneous variation
- dsQTL: genetics of chromatin accessibility – unraveling eQTL mechanism?
- CCLE: reproducing an application of elasticnet (combining lasso and ridge regression) to tumor chemosensitivity



**Figure 1.** Plausible sites of action for genetic determinants of mRNA levels. Genetic variations influencing gene expression may reside within the regulatory sequences, promoters, enhancers, splice sites, and secondary structure motifs of the target gene and so be genetically in *cis* (red stars), or there may be variations in the molecular machinery that interact with *cis*-regulatory sequences and so act genetically in *trans* (blue stars).

# Opportunities for greedy tuning of cis-eQTL search



---

OPINION

# Tackling the widespread and critical impact of batch effects in high-throughput data

---

*Jeffrey T. Leek, Robert B. Scharpf, Héctor Corrada Bravo, David Simcha, Benjamin Langmead, W. Evan Johnson, Donald Geman, Keith Baggerly and Rafael A. Irizarry*

**Abstract** | High-throughput technologies are widely used, for example to assay genetic variants, gene and protein expression, and epigenetic modifications. One often overlooked complication with such studies is batch effects, which occur because measurements are affected by laboratory conditions, reagent lots and personnel differences. This becomes a major problem when batch effects are correlated with an outcome of interest and lead to incorrect conclusions. Using both published studies and our own analyses, we argue that batch effects (as well as other technical and biological artefacts) are widespread and critical to address. We review experimental and computational approaches for doing so.

Many technologies used in biology — and hardware, along with highly trained per-

affected by both biological and non-biological factors. Here we focus on batch effects, a common and powerful source of variation in high-throughput experiments.

Batch effects are sub-groups of measurements that have qualitatively different behaviour across conditions and are unrelated to the biological or scientific variables in a study. For example, batch effects may occur if a subset of experiments was run on Monday and another set on Tuesday, if two technicians were responsible for different subsets of the experiments or if two different lots of reagents, chips or instruments were used. These effects are not exclusive to high-throughput biology and genomics research<sup>1</sup>, and batch effects also affect low-dimensional molecular measurements, such as northern blots and quantitative PCR. Although batch effects are difficult or impossible to detect in low-dimensional assays, high-throughput technologies provide enough data to detect and even remove them. However, if not properly dealt with, these effects can have a particularly strong and pervasive impact. Specific examples have been documented in published studies<sup>2,3</sup> in which the biologi-

Table 1 | **Batch effects seen for a range of high-throughput technologies**

Study description*	Known variable used as a surrogate			Principal components used as a surrogate			Association with outcome Significant features (%) <sup>††</sup>	Refs
	Surrogate <sup>‡</sup>	Confounding (%) <sup>§</sup>	Susceptible features (%) <sup>  </sup>	Principal components rank of surrogate (correlation) <sup>¶</sup>	Principal components rank of outcome (correlation) <sup>*</sup>	Susceptible features (%) <sup>**</sup>		
Data set 1: gene expression microarray, Affymetrix ( $N_p = 22,283$ )	Date	29.7	50.5	1 (0.570)	1 (0.649)	91.6	71.9	9
Data set 2: gene expression, Affymetrix ( $N_p = 4167$ )	Date	77.6	73.7	1 (0.922)	1 (0.668)	98.5	62.2	2
Data set 3: mass spectrometry ( $N_p = 15,154$ )	Processing group	100	51.7	2 (0.344)	2 (0.344)	99.7	51.7	3
Data set 4: copy number variation, Affymetrix ( $N_p = 945,806$ )	Date	29.2	99.5	2 (0.921)	3 (0.485)	99.8	98.8	16
Data set 5: copy number variation, Affymetrix ( $N_p = 945,806$ )	Date	12.2	83.8	1 (0.553)	1 (0.137)	99.8	74.1	17
Data set 6: gene expression, Affymetrix ( $N_p = 22,277$ )	Processing group	NA	83.8	5 (0.369)	NA	97.1	NA	18
Data set 7: gene expression, Agilent ( $N_p = 17,594$ )	Date	NA	62.8	2 (0.248)	NA	96.7	NA	18

ment will be highly correlated with cancer status. Principal components capture both biological and technical variability and, in some cases, principal components can be estimated after the biological variables have been accounted for<sup>15</sup>. In this case, the principal components primarily quantify the effects of artefacts on the high-throughput data. Principal components can be compared to known variables, such as processing group or time. If the principal components do not correlate with these known variables, there may be an alternative, unmeasured source of batch effects in the data.

Involving the sTCC study, we examined the extent of batch effects for eight other published or publicly available data sets (TABLE 1) using the following approach. First, we identified a surrogate for batch effects (such as date or processing group) for each data set. We then used simple linear models to measure the level of confounding between this surrogate and the study outcome (for example, case or control) when available. Note that the more confounding there is, the more likely it is that batch variability can be confused with biological variability. We then summarized the



# Upshots

- eQTLs are, in principle, identifiable by simple linear modeling of relationship between average expression and SNP genotype
- There is specificity to tissues, ...
- Works of Stegle, Storey, Leek et al. indicate that removal of PCs and allied factors from expression array archives is important for improving sensitivity of eQTL detection

# dsQTL identification

## LETTER

doi:10.1038/nature10808

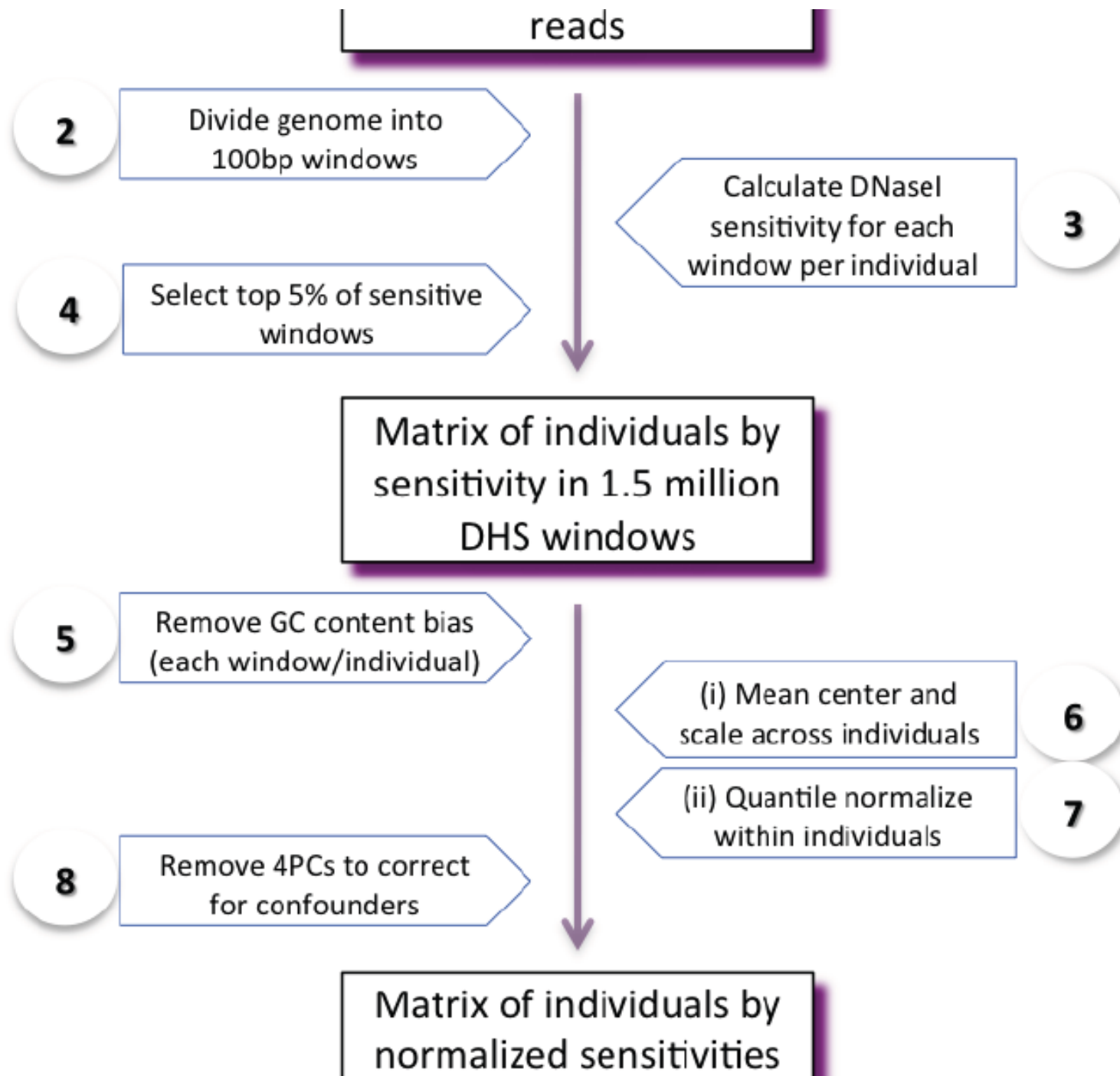
### DNase I sensitivity QTLs are a major determinant of human expression variation

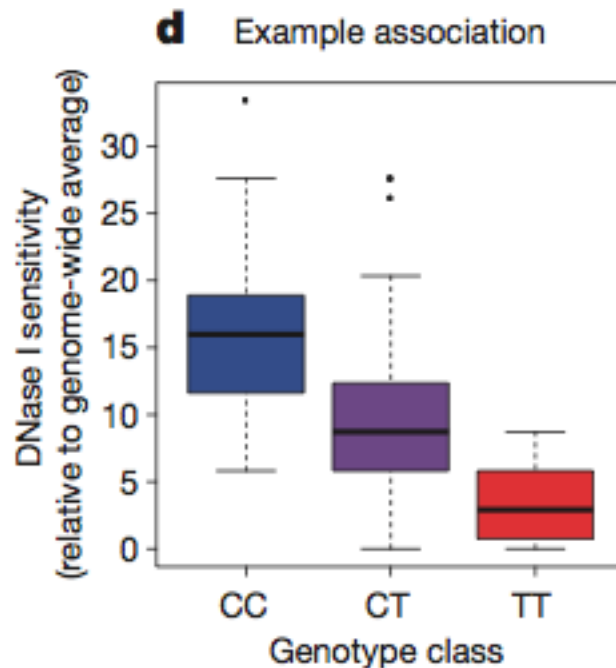
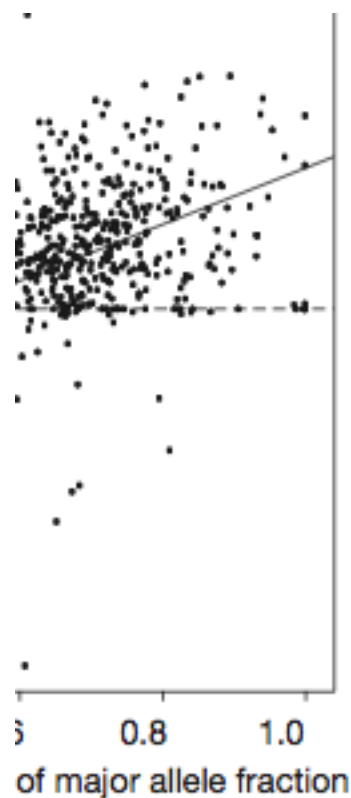
Jacob F. Degner<sup>1,2\*</sup>, Athma A. Pai<sup>1\*</sup>, Roger Pique-Regi<sup>1\*</sup>, Jean-Baptiste Veyrieras<sup>1,3</sup>, Daniel J. Gaffney<sup>1,4</sup>, Joseph K. Pickrell<sup>1</sup>, Sherryl De Leon<sup>4</sup>, Katelyn Michelini<sup>4</sup>, Noah Lewellen<sup>4</sup>, Gregory E. Crawford<sup>5,6</sup>, Matthew Stephens<sup>1,7</sup>, Yoav Gilad<sup>1</sup> & Jonathan K. Pritchard<sup>1,4</sup>

The mapping of expression quantitative trait loci (eQTLs) has emerged as an important tool for linking genetic variation to changes in gene regulation<sup>1-5</sup>. However, it remains difficult to identify the causal variants underlying eQTLs, and little is known about the regulatory mechanisms by which they act. Here we show that genetic variants that modify chromatin accessibility and transcription factor binding are a major mechanism through which genetic variation leads to gene expression differences among humans. We used DNase I sequencing to measure chromatin accessibility in 70 Yoruba lymphoblastoid cell lines, for which genome-wide genotypes and estimates of gene expression levels are also available<sup>6-8</sup>. We obtained a total of 2.7 billion uniquely

and enhancer-associated histone marks. Furthermore, bound transcription factors protect the DNA sequence within a binding site from DNase I cleavage, often producing recognizable 'footprints' of decreased DNase I sensitivity<sup>13,15-17</sup>.

We collected DNase-seq data for 70 HapMap Yoruba lymphoblastoid cell lines for which gene expression data and genome-wide genotypes were already available<sup>6-8</sup>. We obtained an average of 39 million uniquely mapped DNase-seq reads per sample, providing individual maps of chromatin accessibility for each cell line (see Supplementary Information for all analysis details). Our data allowed us to characterize the distribution of DNase I cuts within individual hypersensitive sites at extremely high resolution. As expected, the DHSs coincided to a great





Position relative to centromere



**Association of dsQTLs and a typical example.**  
 Association between DNase I cut rates in 100-bp  
 regions (green) and 40-kb (black) regions centred  
 on the same SNP. Allele-specific analysis of dsQTLs in

dsQTL (rs4953223). The bla  
**d**, Box plot showing that rs4  
 accessibility ( $P = 3 \times 10^{-13}$   
 DNase I sensitivity, disrupts

# Different approaches to dsQTL data representation

- Chicago/GEO
  - Filtered and normalized DHS assay results in 70 bed files, indexed to hg18: 1.4GB gzipped on GEO, metadata not directly bound
  - Imputed genotype data harbored separately as name, loc, alleles, expected B allele count per indiv/SNP: chr1 = .3GB gzipped, 5 text bytes/SNP
- Bioconductor (dsQTLtools, not posted yet)
  - .8GB compressed `SummarizedExperiment` for DHS plus .08GB for 4 million imputed genotypes

# SummarizedExperiment instance

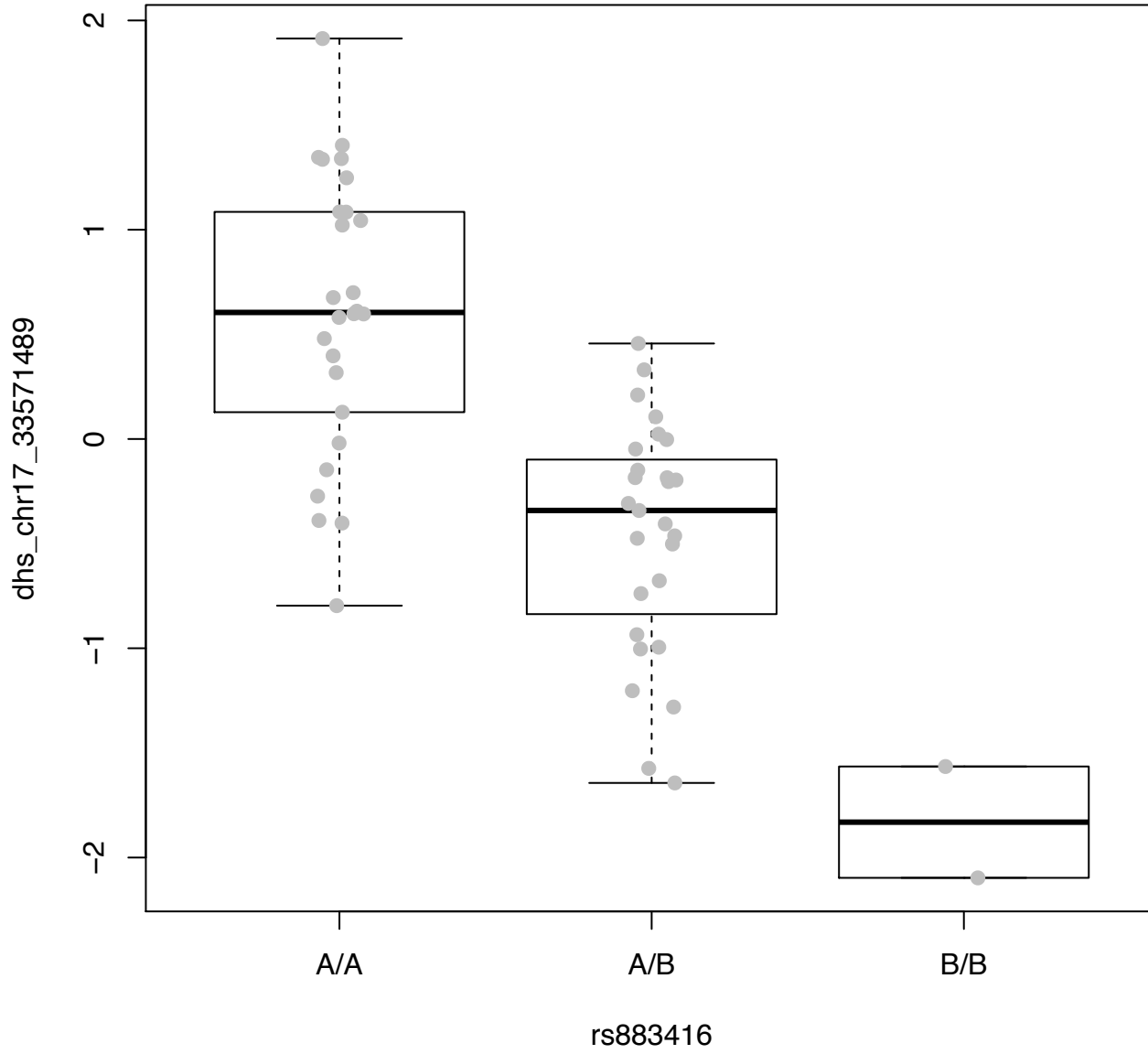
```
> DHStop5_hg19
class: SummarizedExperiment
dim: 1465442 70
exptData(2): MIAME annotation
assays(1): scores
rownames(1465442): dhs_chr1_10402 dhs_chr1_10502 ...
  dhs_chr22_51228236 dhs_chr22_51234736
rowData metadata column names(0):
colnames(70): NA18486 NA18498 ... NA19239 NA19257
colData names(9): naid one ... male isFounder
```

```
> assays(DHStop5_hg19)$scores[1:5,1:3]
              NA18486      NA18498      NA18499
dhs_chr1_10402 -0.8932210 -0.3633581 -0.4540041
dhs_chr1_10502 -0.1523477 -0.1704101 -1.0598971
dhs_chr1_13239  0.4360728 -0.1159094  1.2505193
dhs_chr1_13939 -0.5259945 -0.8212344  0.1145535
dhs_chr1_16039 -0.9991160  0.2092481  0.3199874
```

```
> s1 = dsqNearGene("SLFN5")
> s1
dsqLook instance for SLFN5 w/ radius 1000.
best DHS site: dhs_chr17_33571489.
```

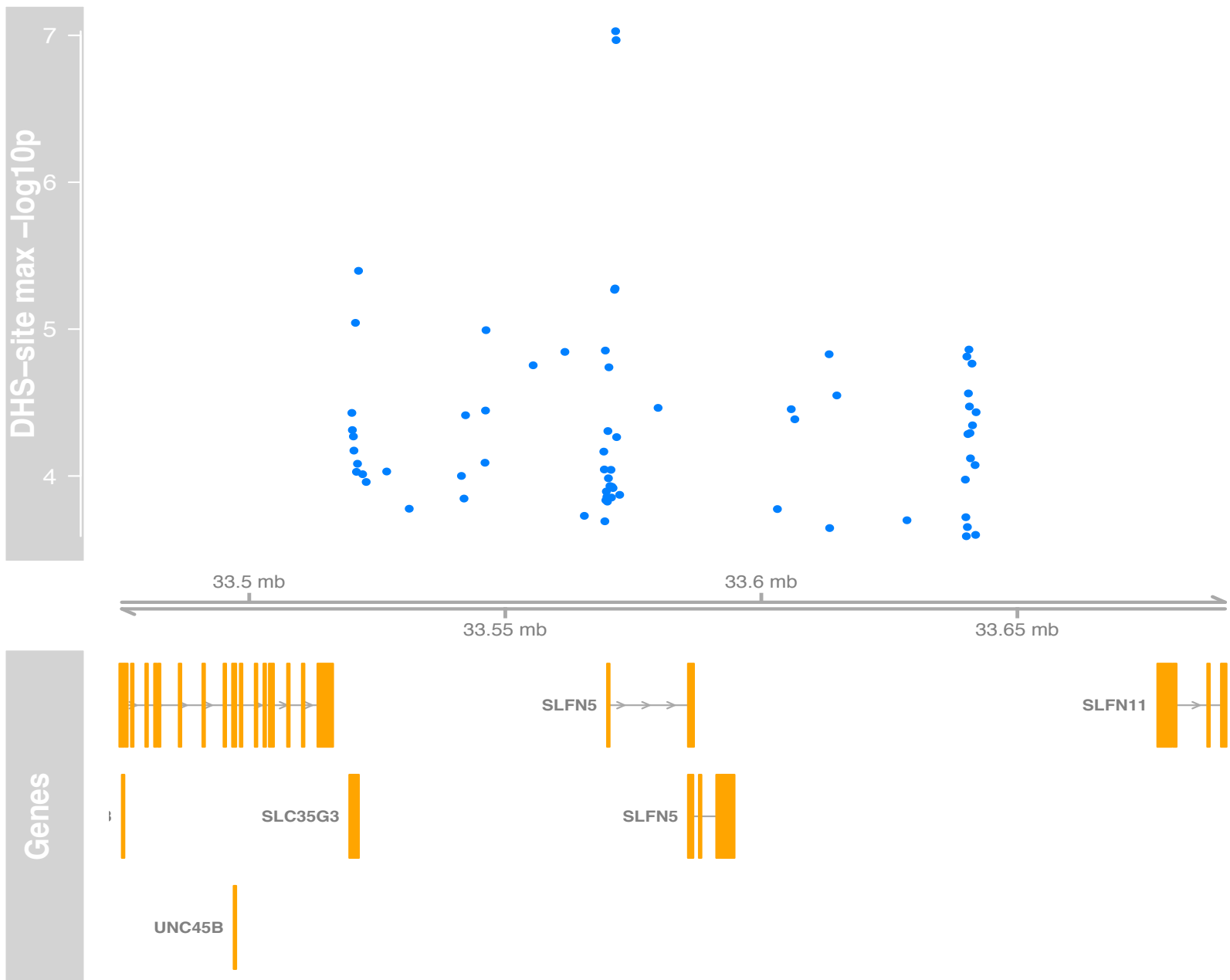
- R Under development (unstable) (2013-01-08 r61589),  
x86\_64-apple-darwin10.8.0
- Locale:  
en\_US.US-ASCII/en\_US.US-ASCII/en\_US.US-ASCII/C/en\_US.US-ASCII/en\_US.US-ASCII
- Base packages: base, datasets, graphics, grDevices, methods, parallel, splines,  
stats, stats4, tools, utils
- Other packages: AnnotationDbi 1.21.9, Biobase 2.19.2, BiocGenerics 0.5.6,  
BiocInstaller 1.9.6, Biostrings 2.27.8, codetools 0.2-8, DBI 0.2-5, digest 0.6.0,  
dsQTLtools 0.0.5, GenomicFeatures 1.11.6, GenomicRanges 1.11.21,  
GGBase 3.21.2, GGtools 4.7.17, GO.db 2.8.0, hmyriB36 0.99.16,  
Homo.sapiens 1.0.0, IRanges 1.17.24, lattice 0.20-13, Matrix 1.0-10,  
org.Hs.eg.db 2.8.0, OrganismDbi 1.1.9, Rsamtools 1.11.14, RSQlite 0.11.2,  
snpStats 1.9.2, survival 2.37-2, TxDb.Hsapiens.UCSC.hg19.knownGene 2.8.0,  
weaver 1.25.0
- Loaded via a namespace (and not attached): annotate 1.37.3, biomaRt 2.15.0,  
bit 1.1-9, bitops 1.0-5, BSgenome 1.27.1, ff 2.2-10, genefilter 1.41.1, graph 1.37.4,  
grid 3.0.0, RBGL 1.35.0, RCurl 1.95-3, rtracklayer 1.19.6,  
VariantAnnotation 1.5.28, XML 3.95-0.1, xtable 1.7-0, zlibbioc 1.5.0

best dsQTL near SLFN5 rad. 50000  
-log10 assoc p. = 7.028





# dsQTL scores



# AnnotationHub for DnaseI peaks

```
> library(AnnotationHub)
> ah = AnnotationHub()
> nah = names(ah)
> ds = grep("dnase", nah, ignore.case=TRUE, value=TRUE)
> length(ds)
[1] 686
> ds[1:6]
[1] "goldenpath.dm2.database.bdtnpDnase_0.0.1.RData"
[2] "goldenpath.hg16.database.nhgriDnaseHs_0.0.1.RData"
[3] "goldenpath.hg17.database.sangamoDnaseHs_0.0.1.RData"
[4] "goldenpath.hg18.database.wgEncodeUwDnaseSeq_0.0.1.RData"
[5] "goldenpath.hg19.encodeDCC.wgEncodeAvgDnaseUniform.wgEncodeAvgDnaseDuke8988t
UniPk.narrowPeak_0.0.1.RData"
[6] "goldenpath.hg19.encodeDCC.wgEncodeAvgDnaseUniform.wgEncodeAvgDnaseDukeAosmc
UniPk.narrowPeak_0.0.1.RData"
> ds[41:45]
[41] "goldenpath.hg19.encodeDCC.wgEncodeAvgDnaseUniform.wgEncodeAvgDnaseDuke8988t
UniPk.narrowPeak_0.0.1.RData"
[42] "goldenpath.hg19.encodeDCC.wgEncodeAvgDnaseUniform.wgEncodeAvgDnaseDukeAosmc
UniPk.narrowPeak_0.0.1.RData"
[43] "goldenpath.hg19.encodeDCC.wgEncodeAvgDnaseUniform.wgEncodeAvgDnaseDuke8988t
UniPk.narrowPeak_0.0.1.RData"
[44] "goldenpath.hg19.encodeDCC.wgEncodeAvgDnaseUniform.wgEncodeAvgDnaseDukeAosmc
UniPk.narrowPeak_0.0.1.RData"
[45] "goldenpath.hg19.encodeDCC.wgEncodeAvgDnaseUniform.wgEncodeAvgDnaseDuke8988t
UniPk.narrowPeak_0.0.1.RData"
```

# Current metadata

```
> dst1 = ah[[ds[5]]]
Retrieving 'goldenpath/hg19/encodeDCC/wgEncodeAvgDnaseUniform/wgEncodeAvgDnaseDuke8988tUniPk.narrowPeak_0.0.1.RData'
> args(ahinfo)
function (hub, path)
NULL
> ahinfo(ah, ds[5] )
From: EncodeDCC
Version: ENCODE Jan 2011 Freeze
Description: wgEncodeAvgDnaseDuke8988tUniPk
Genus and Species: Homo sapiens
Genome: hg19
BiocVersion: 2.12, 2.13
Tags: 8988T, wgEncodeAvgDnaseUniPk, DnaseSeq, ENCODE Jan 2011 Freeze, wgEncodeE
H001103, Duke, 80fadeb7a14a72add38203910d937f50, wgEncode, 1700000, wgEncodeAvgD
naseDuke8988tUniPk, None, narrowPeak, Peaks, wgEncodeAvgDnaseUniform
```

# Upshots

- Representation of Dnase1 HS can take various forms at various scales, tracks are nice but archive is complicated
- DS-seq archive very substantial even as filtered, but a SummarizedExperiment container can manage it
- Searching for dsQTL with substantial parallelism and small RAM footprint: Martin and Val's Streamer, scanVCF

# Some machine learning with CCLE

## LETTER

doi:10.1038/nature11003

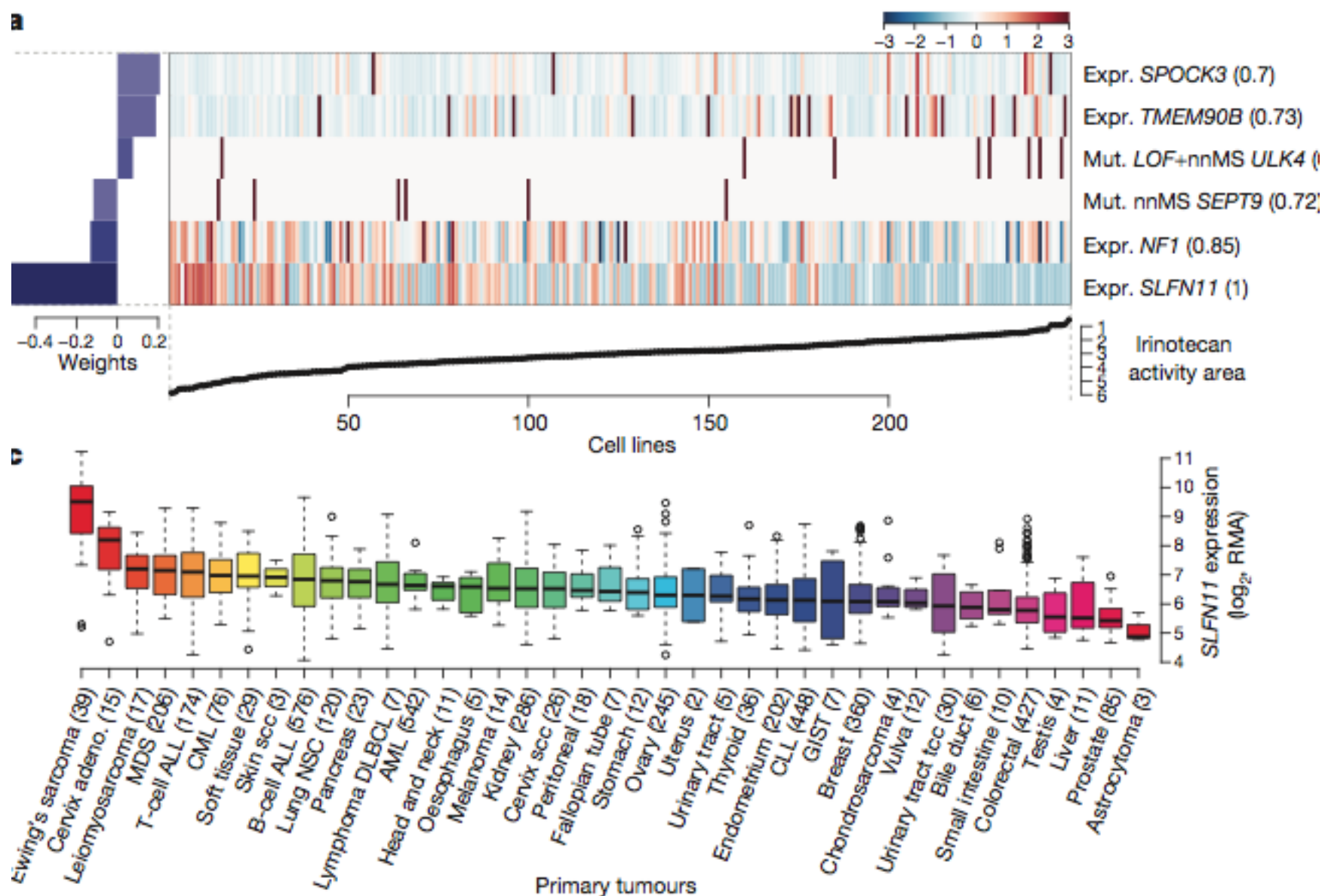
### The Cancer Cell Line Encyclopedia enables predictive modelling of anticancer drug sensitivity

Jordi Barretina<sup>1,2,3†\*</sup>, Giordano Caponigro<sup>4\*</sup>, Nicolas Stransky<sup>1\*</sup>, Kavitha Venkatesan<sup>4\*</sup>, Adam A. Margolin<sup>1†\*</sup>, Sungjoon Kim<sup>5</sup>, Christopher J. Wilson<sup>4</sup>, Joseph Lehár<sup>4</sup>, Gregory V. Kryukov<sup>1</sup>, Dmitriy Sonkin<sup>4</sup>, Anupama Reddy<sup>4</sup>, Manway Liu<sup>4</sup>, Lauren Murray<sup>1</sup>, Michael F. Berger<sup>1†</sup>, John E. Monahan<sup>4</sup>, Paula Morais<sup>1</sup>, Jodi Meltzer<sup>4</sup>, Adam Korejwa<sup>1</sup>, Judit Jané-Valbuena<sup>1,2</sup>, Felipa A. Mapa<sup>4</sup>, Joseph Thibault<sup>5</sup>, Eva Bric-Furlong<sup>4</sup>, Pichai Raman<sup>4</sup>, Aaron Shipway<sup>5</sup>, Ingo H. Engels<sup>5</sup>, Jill Cheng<sup>6</sup>, Guoying K. Yu<sup>6</sup>, Jianjun Yu<sup>6</sup>, Peter Aspesi Jr<sup>4</sup>, Melanie de Silva<sup>4</sup>, Kalpana Jagtap<sup>4</sup>, Michael D. Jones<sup>4</sup>, Li Wang<sup>4</sup>, Charles Hatton<sup>3</sup>, Emanuele Palescandolo<sup>3</sup>, Supriya Gupta<sup>1</sup>, Scott Mahan<sup>1</sup>, Carrie Sougnez<sup>1</sup>, Robert C. Onofrio<sup>1</sup>, Ted Liefeld<sup>1</sup>, Laura MacConaill<sup>3</sup>, Wendy Winckler<sup>1</sup>, Michael Reich<sup>1</sup>, Nanxin Li<sup>5</sup>, Jill P. Mesirov<sup>1</sup>, Stacey B. Gabriel<sup>1</sup>, Gad Getz<sup>1</sup>, Kristin Ardlie<sup>1</sup>, Vivien Chan<sup>6</sup>, Vic E. Myer<sup>4</sup>, Barbara L. Weber<sup>4</sup>, Jeff Porter<sup>4</sup>, Markus Warmuth<sup>4</sup>, Peter Finan<sup>4</sup>, Jennifer L. Harris<sup>5</sup>, Matthew Meyerson<sup>1,2,3</sup>, Todd R. Golub<sup>1,3,7,8</sup>, Michael P. Morrissey<sup>4\*</sup>, William R. Sellers<sup>4\*</sup>, Robert Schlegel<sup>4\*</sup> & Levi A. Garraway<sup>1,2,3\*</sup>

The systematic translation of cancer genomic data into knowledge of tumour biology and therapeutic possibilities remains challenging. Such efforts should be greatly aided by robust preclinical model systems that reflect the genomic diversity of human cancers and for which detailed genetic and pharmacological annotation is available<sup>1</sup>.

Here we describe the Cancer Cell Line Encyclopedia (CCLE): a

known cancer genes were assessed by mass spectrometric genotyping<sup>13</sup> (Supplementary Table 2 and Supplementary Fig. 1). DNA copy number was measured using high-density single nucleotide polymorphism arrays (Affymetrix SNP 6.0; Supplementary Methods). Finally, messenger RNA expression levels were obtained for each of the lines using Affymetrix U133 plus 2.0 arrays. These data were also used to confirm cell line



**Figure 4 | Predicting sensitivity to topoisomerase I inhibitors.** **a**, Elastic net regression analysis of genomic correlates of irinotecan sensitivity is shown for 250 cell lines. **b**, Dose–response curves for three Ewing’s sarcoma cell lines

the mean growth inhibition ( $n$  tumours). Box-and-whisker plot each subtype, ordered by the  $n$

## Specification [\[edit\]](#)

---

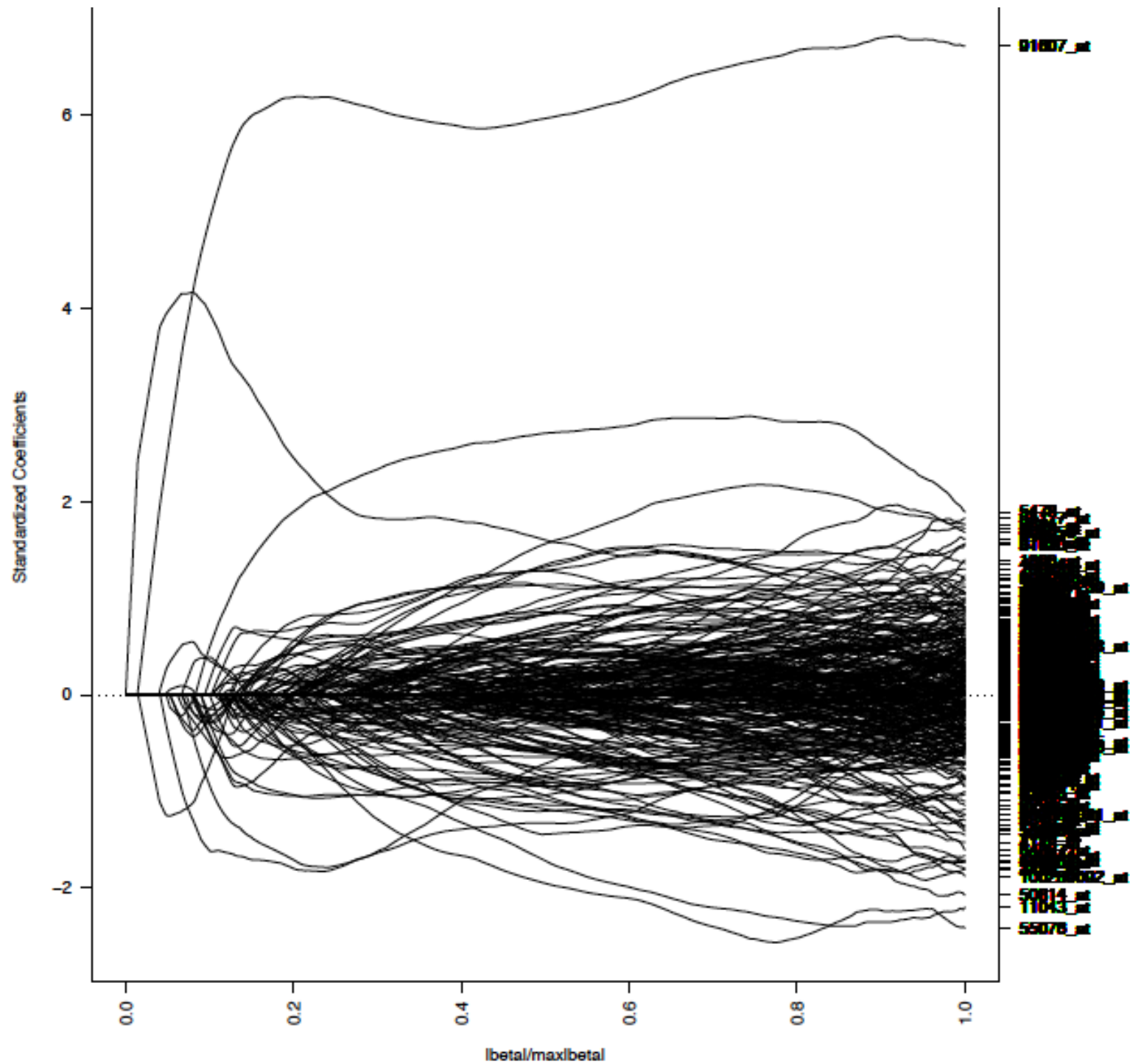
The elastic net method which overcomes the limitations of the [LASSO \(least absolute shrinkage and selection operator\)](#) method which uses a penalty function based on

$$\|\beta\|_1 = \sum_{j=1}^p |\beta_j|.$$

Use of this penalty function has several limitations.<sup>[1]</sup> For example, in the "large  $p$ , small  $n$  problem" case, the LASSO selects at most  $n$  variables before it saturates. Also if there is a group of highly correlated variables, then the LASSO tends to select one variable from a group and ignore the others. To overcome these limitations, the elastic net adds a quadratic part to the penalty ( $\|\beta\|^2$ ), which when used alone is [ridge regression](#) (known also as Tikhonov regularization). The estimates from the elastic net method are defined by

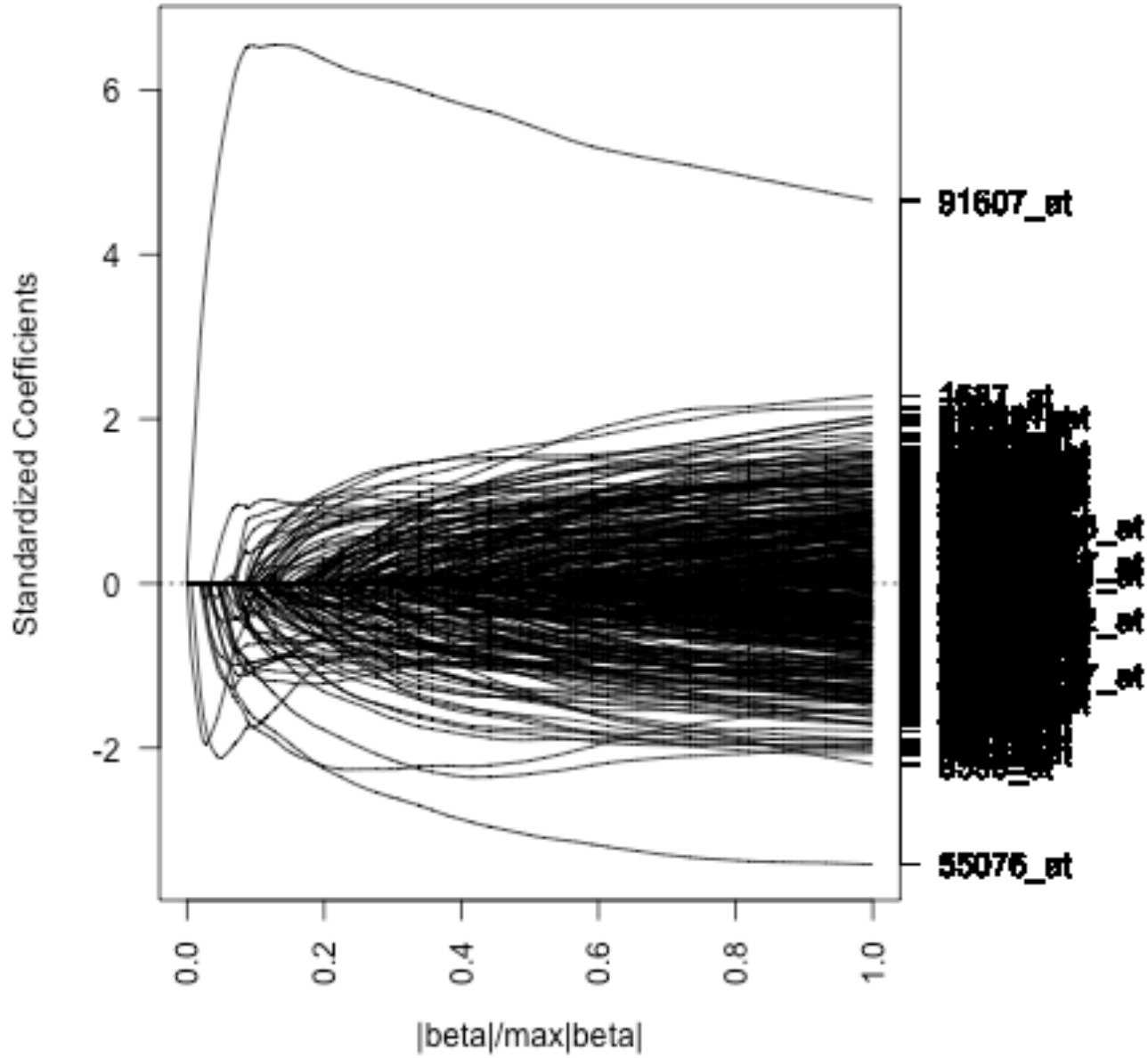
$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}}(\|y - X\beta\|^2 + \lambda_2\|\beta\|^2 + \lambda_1\|\beta\|_1).$$

As a result, the elastic net method includes the LASSO and ridge regression: in other words, each of them is a special case where  $\lambda_1 = 1, \lambda_2 = 0$  or  $\lambda_1 = 0, \lambda_2 = 1$ . Meanwhile, the naive version of elastic net method finds an estimator in a two-stage procedure : first for each fixed  $\lambda_2$  it finds the ridge regression coefficients, and then does a LASSO type shrinkage. This kind of estimation incurs a double amount of shrinkage, which introduces unnecessary extra bias and outcomes with bad prediction performance. To improve the prediction performance, the authors rescale the coefficients of the naive version of elastic net by multiplying the estimated coefficients by  $(1 + \lambda_2)$ .<sup>[1]</sup>

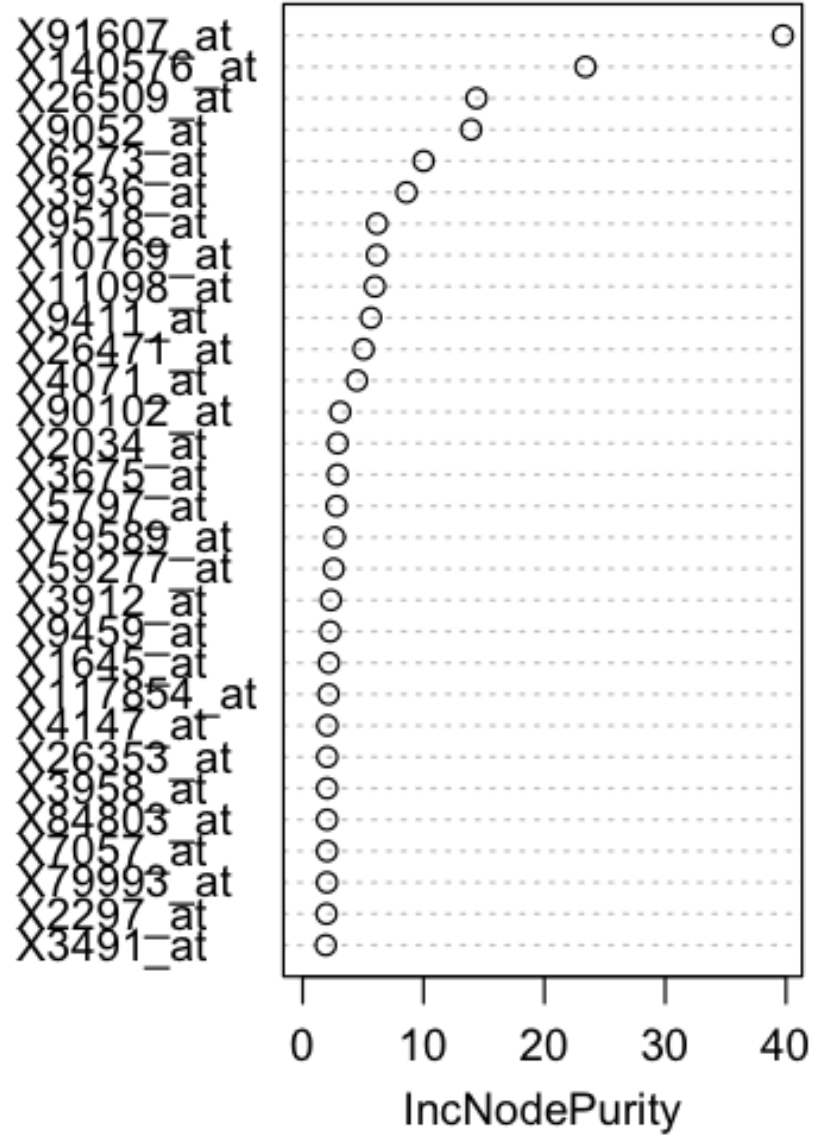
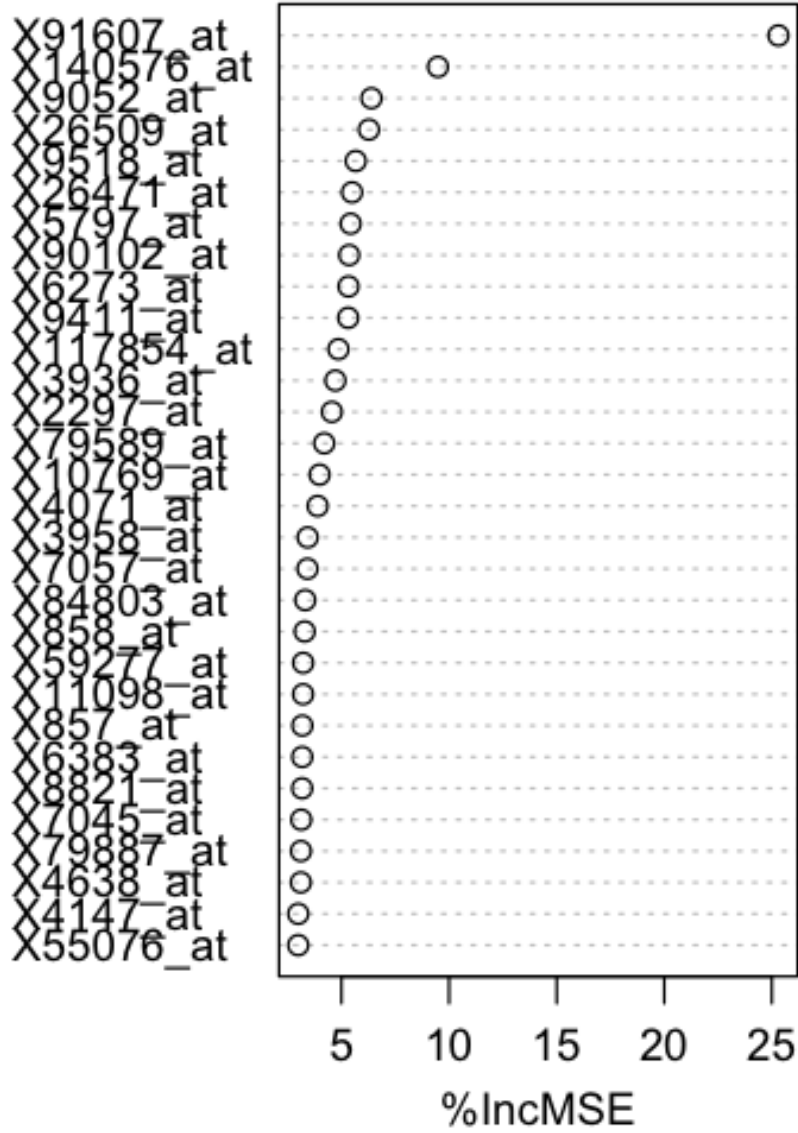




lam = .1



rf1



# Comments

- Qualitative effects of tuning the elastic net vs. lasso ... organize with CV ... performance of a single fit not so great
- Please consider models beyond “main effects”  
– randomForests is a black-box approach that accommodates one approach to variable interaction
- What about batch effects? Has the published expression data been properly adjusted?

# Analysis Tools: Gene Set Enrichment Analysis (GSEA)

[PAGE INFO ▲](#)

**Step1: Tools** | **Step2: Sample Set** | **Step3: Create Sample Set**

## Step3: Create a New Sample Set

**STEP INFO ▲**

**First class for comparative analyses - [Close](#)**

**Class 1 [edit](#)** **Filter by Mu**

**Filter by:**  Cell line primary name  Cell line aliases  Gender  Site Primary  Histology  Hist Subtype1  Source  Hybrid Capture Se

[\(show\)](#) **1036 Items**

**Second class for comparative analyses - [Close](#)**

**Class 2 [edit](#)** **Filter by Mu**

**Filter by:**  Cell line primary name  Cell line aliases  Gender  Site Primary  Histology  Hist Subtype1  Source  Hybrid Capture Se

[\(show\)](#) **1036 Items**

# Quiz

- How should we organize data from integrative experiments (expression, CNV, genotype, drug sensitivity)?
- Ordering genes measured across tumor types with respect to association between expression and drug sensitivity: what methods are preferred? Can we distinguish sensitive from insensitive tumor types?
- How to test whether a given mutation distinguishes sensitivities within a tumor class?