# GS01 0163
# Analysis of Microarray Data

Keith Baggerly and Kevin Coombes
Section of Bioinformatics
Department of Biostatistics and Applied Mathematics
UT M. D. Anderson Cancer Center

kabagg@mdanderson.org

kcoombes@mdanderson.org

18 October 2005

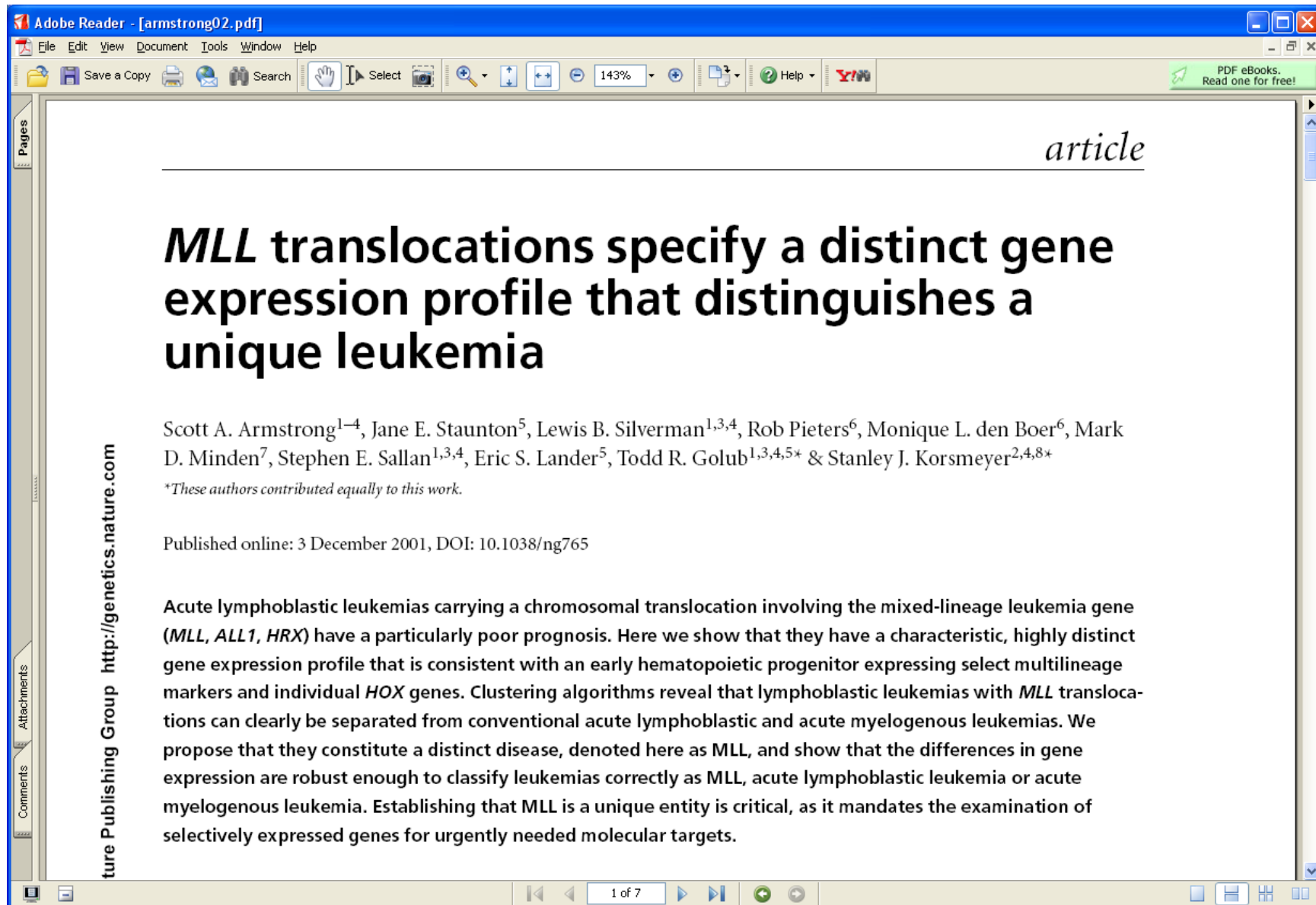# Lecture 13: Microarrays in R: Start to Finish

- Source of the Data Set

- Understanding the Sample Information

- First Pass Using dChip

- Starting in R

- Just RMA

# Source of the Data Set

In today's lecture, we're going to perform a complete start-to-finish analysis of a microarray data set. We have chosen to use a leukemia data set that we looked at briefly in an earlier lecture. The data set consists of U95A microarray experiments on

1. 24 patients with acute lymphocytic leukemia (ALL)

2. 28 patients with acute myeloid leukemia (AML)

3. 20 patients with mixed lineage leukemia (MLL)

# Armstrong, Nat Genet, 2002; 30:41-47

*article*

## MLL translocations specify a distinct gene expression profile that distinguishes a unique leukemia

Scott A. Armstrong[1–4], Jane E. Staunton[5], Lewis B. Silverman[1,3,4], Rob Pieters[6], Monique L. den Boer[6], Mark D. Minden[7], Stephen E. Sallan[1,3,4], Eric S. Lander[5], Todd R. Golub[1,3,4,5*] & Stanley J. Korsmeyer[2,4,8*]

*These authors contributed equally to this work.*

Acute lymphoblastic leukemias carrying a chromosomal translocation involving the mixed-lineage leukemia gene (MLL, ALL1, HRX) have a particularly poor prognosis. Here we show that they have a characteristic, highly distinct gene expression profile that is consistent with an early hematopoietic progenitor expressing select multilineage markers and individual HOX genes. Clustering algorithms reveal that lymphoblastic leukemias with MLL translocations can clearly be separated from conventional acute lymphoblastic and acute myelogenous leukemias. We propose that they constitute a distinct disease, denoted here as MLL, and show that the differences in gene expression are robust enough to classify leukemias correctly as MLL, acute lymphoblastic leukemia or acute myelogenous leukemia. Establishing that MLL is a unique entity is critical, as it mandates the examination of selectively expressed genes for urgently needed molecular targets.

ture Publishing Group  http://genetics.nature.com

# Results reported in the paper

- MLL is distinct from conventional ALL

  - Based on arrays, ~1200 differentially expressed genes

- MLL shows multilineage gene expression

  - Looking at the expressed genes, find some from B cells and some from myeloid cells

- MLL is arrested at an early stage of hematopoiesis

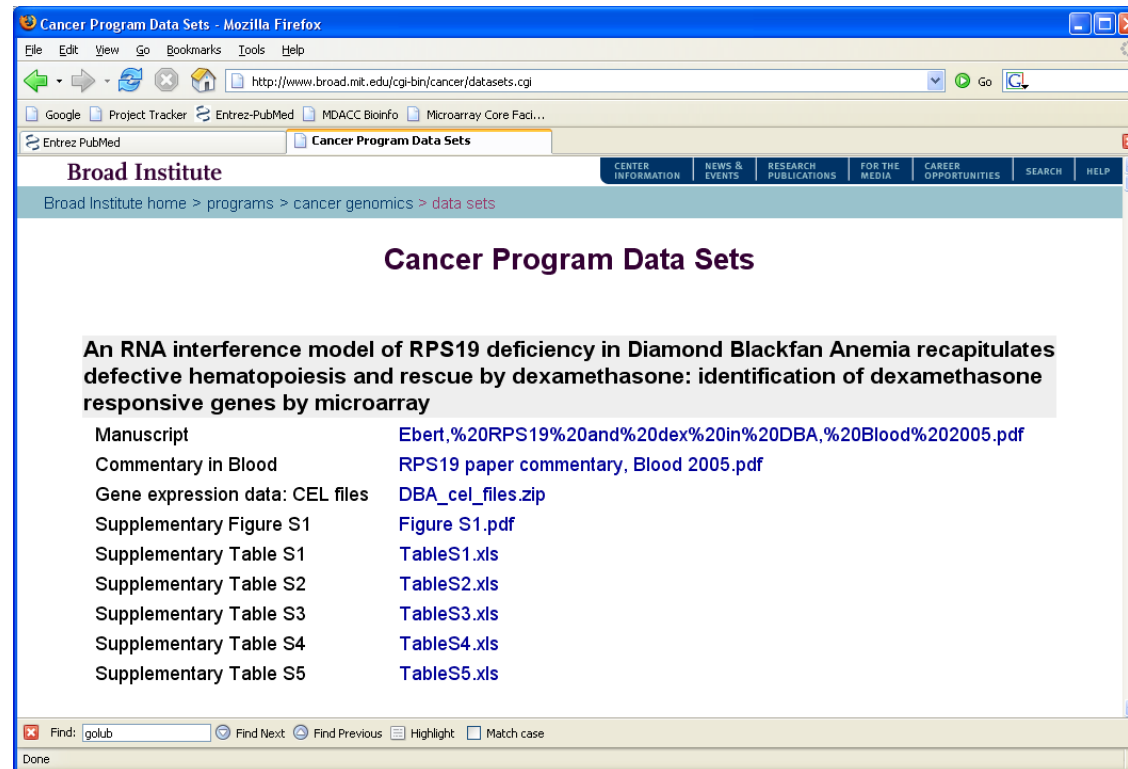- Some HOX genes are overexpressed in MLL

# Results reported in the paper

- MLL is distinct from AML as well as from conventional ALL

    - Principal components analysis
    - Selected genes using one-versus-all comparisons

- Gene expression profiles correctly classify ALL, AML, MLL

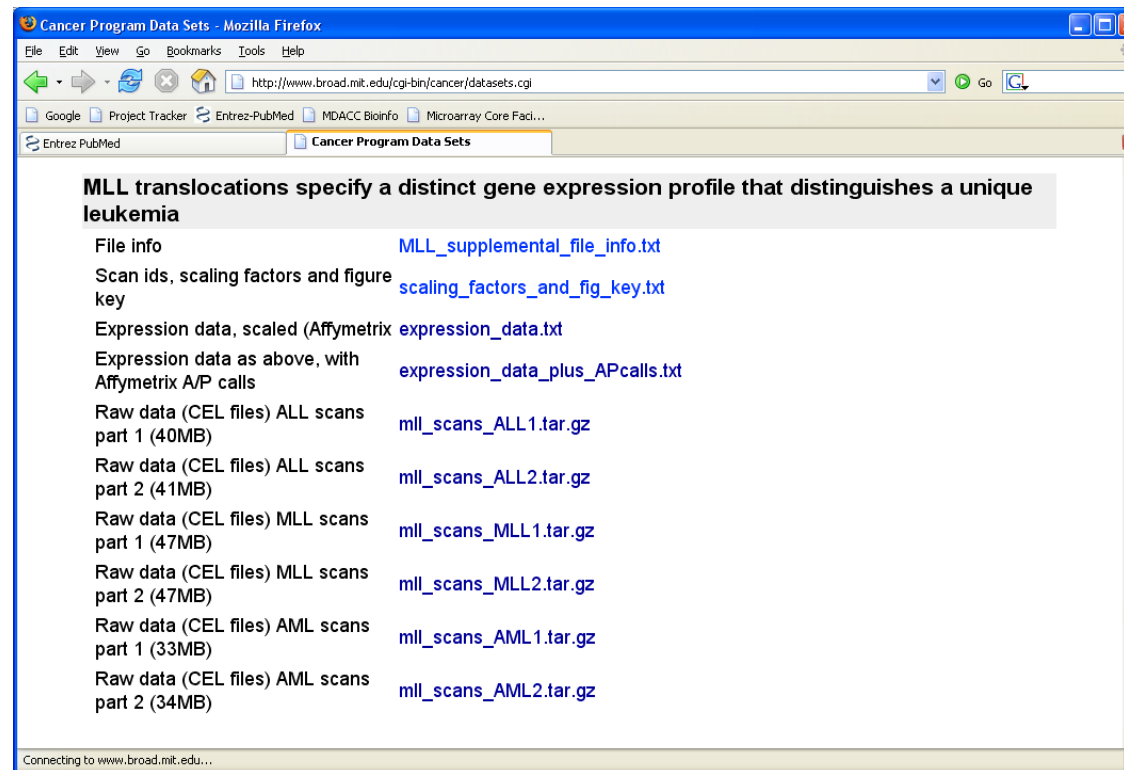    - K-nearest-neighbors (KNN) predictor

# Getting the data

Note: The links to supplementary data listed in the original paper no longer work (as of the morning of 18 October 2005). However, the data is available at: `http:`
`//www.broad.mit.edu/cgi-bin/cancer/datasets.cgi`

# Getting the data

Scroll down (or search for "translocations") to find the data set. You'll need all six collections of CEL files along with the "scaling_factors_and_fig_key.txt" that contains the sample information.
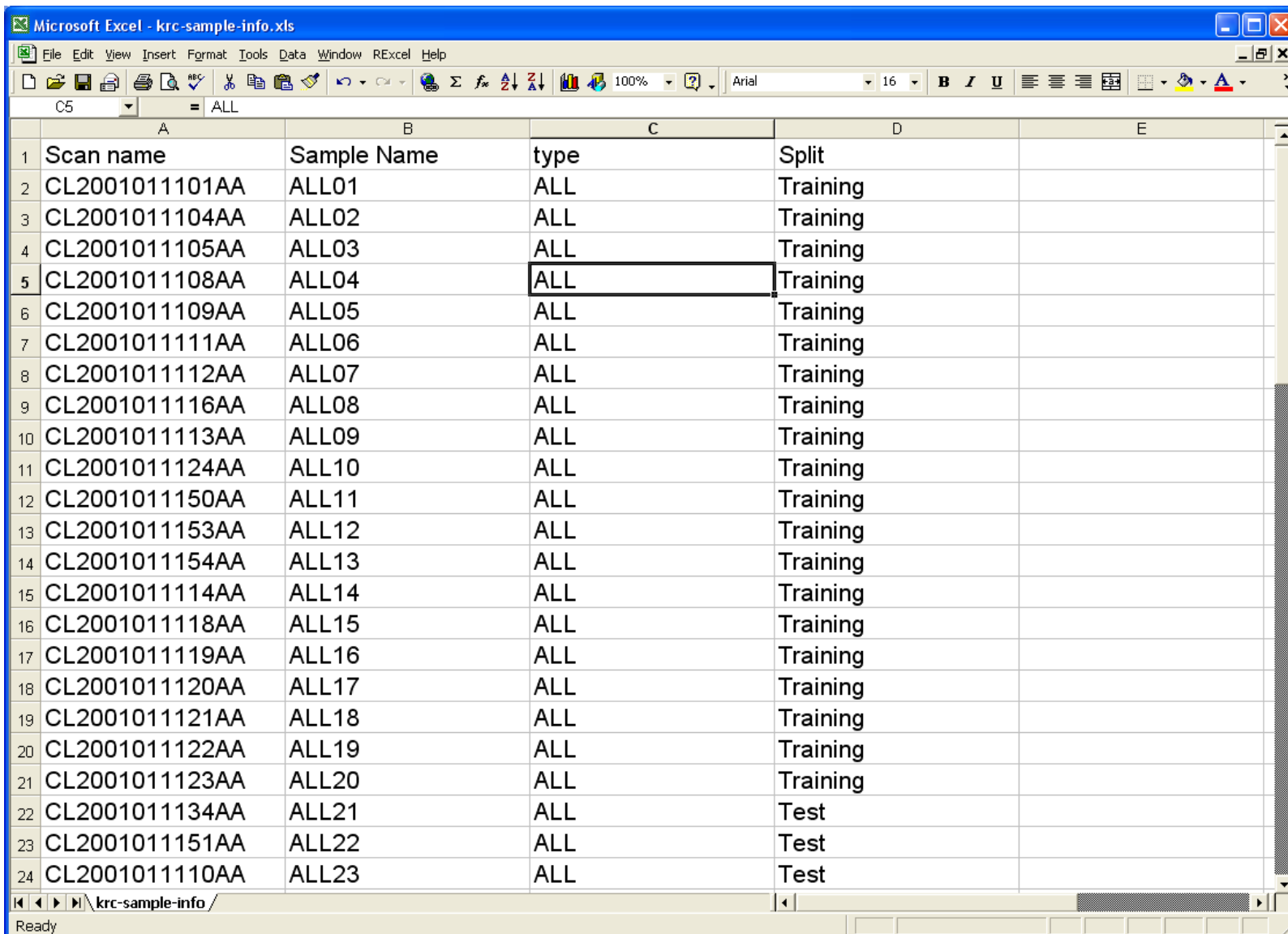
# Understanding the Sample Information

# Including critical factors

# First Pass Using dChip

Early in the course, we used this data set in dChip. One of the things we discovered was that the samples were run on two different iterations of the U95A GeneChip. The ALL and MLL samples were run on version 1 of the U95A and the AML samples were run on the U95Av2 chip.

Our first step in dealing with this issue was to install the CEL files into two different directories. In dChip, we handle this by first producing a "data file list" that tells dChip which directories to look into:



---

GS01 0163: ANALYSIS OF MICROARRAY DATA

# Probe Set Mask File

We also have top produce a "probe set mask" file that tells dChip which probe sets to ignore (since they changed from one version of the chip to the next).

```
hg_u95av2 probe set mask.txt - Notepad
File  Edit  Format  View  Help
160030_at
160039_at
160025_at
160024_at
160023_at
160022_at
160021_r_at
160020_at
160037_at
160038_s_at
160043_at
160035_at
160042_s_at
160027_s_at
160041_at
160040_at
160032_at
160033_s_at
160036_at
160034_s_at
160031_at
160029_at
160028_s_at
160026_at
160044_g_at
```

# dChip analysis

1. Load the files into dChip using

   - The sample information file
   - The data file list
   - The U95Av2 CDF file
   - The probe set mask

2. Normalize to array ALL21 (median brightness)

3. Quantify using the PM-only model

# dChip quality check

- dChip flags several arrays as outliers:

  - ALL10, which has the highest overall brightness
  - MLL09, which has the lowest overall broightness
  - ALL04, ALL16, AML08

# dChip Comparison: ALL vs MLL

# 628 differentially expressed genes



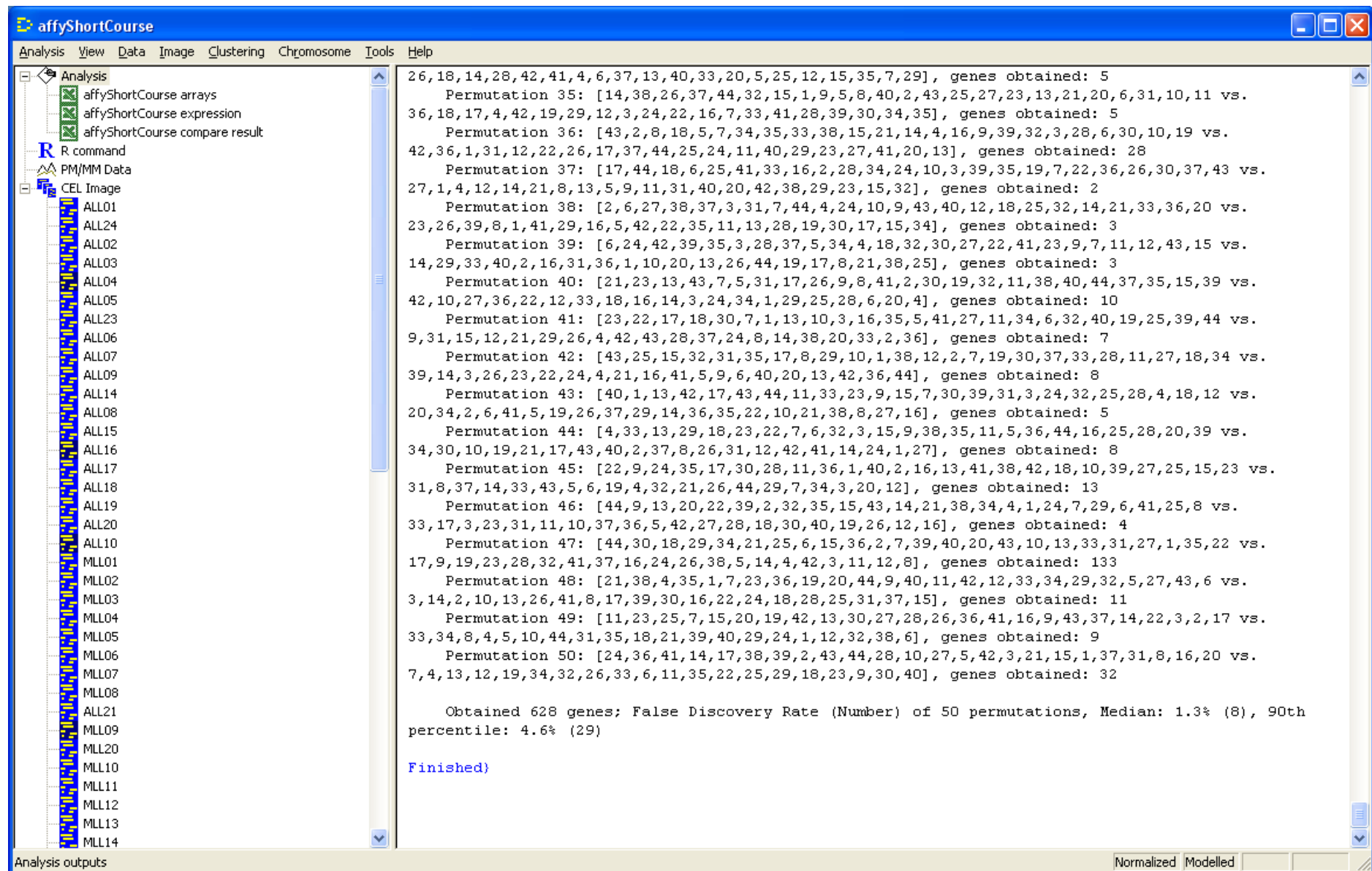| | probe set | gene | baseline mean | experiment mean | fold change | lower bound | upper bound | difference |
|---|---|---|---|---|---|---|---|---|
| 13 | | | | | | | | |
| 14 | 37680_at | A kinase (PRKA) a | 2991.81 | 147.65 | -20.26 | -13.19 | -29.88 | -2844.16 |
| 15 | 37280_at | MAD, mothers agai | 9355.77 | 697.15 | -13.42 | -9.62 | -17.36 | -8658.62 |
| 16 | 1325_at | MAD, mothers agai | 7714.68 | 617.24 | -12.50 | -8.54 | -17.12 | -7097.44 |
| 17 | 1488_at | protein tyrosine ph | 4121.6 | 579.89 | -7.11 | -3.86 | -10.50 | -3541.71 |
| 18 | 34194_at | Homo sapiens mRN | 1420.54 | 200.61 | -7.08 | -3.32 | -13.62 | -1219.93 |
| 19 | 1077_at | recombination activ | 6897.9 | 982.66 | -7.02 | -4.26 | -11.30 | -5915.24 |
| 20 | 753_at | nidogen 2 (osteoni | 2389.96 | 342.13 | -6.99 | -2.78 | -11.31 | -2047.83 |
| 21 | 37908_at | guanine nucleotide | 6909.7 | 1035.23 | -6.67 | -4.15 | -11.67 | -5874.47 |
| 22 | 35614_at | transcription factor | 7521.88 | 1216.05 | -6.19 | -4.07 | -8.59 | -6305.83 |
| 23 | 34800_at | leucine-rich repeat | 5226.23 | 844.98 | -6.19 | -4.00 | -9.80 | -4381.25 |
| 24 | 31892_at | protein tyrosine ph | 868.91 | 145.9 | -5.96 | -2.01 | -10.04 | -723.00 |
| 25 | 41266_at | integrin, alpha 6 | 7952.9 | 1466.32 | -5.42 | -3.85 | -7.47 | -6486.57 |
| 26 | 31786_at | KH domain contain | 2488.25 | 472.29 | -5.27 | -3.29 | -9.34 | -2015.95 |
| 27 | 38408_at | transmembrane 4 s | 6455.32 | 1264.74 | -5.10 | -3.47 | -7.59 | -5190.58 |
| 28 | 36650_at | cyclin D2 | 10001.23 | 2021.61 | -4.95 | -3.47 | -7.76 | -7979.62 |
| 29 | 38578_at | tumor necrosis fact | 4050.73 | 828.6 | -4.89 | -3.33 | -7.03 | -3222.13 |
| 30 | 32778_at | inositol 1,4,5-tripho | 2099.86 | 433.92 | -4.84 | -3.61 | -6.41 | -1665.94 |
| 31 | 40570_at | forkhead box O1A ( | 10307.45 | 2130.08 | -4.84 | -3.37 | -7.52 | -8177.37 |
| 32 | 39878_at | protocadherin 9 | 12520.62 | 2670.21 | -4.69 | -3.07 | -7.31 | -9850.42 |
| 33 | 31886_at | 5'-nucleotidase, ec | 2593.49 | 567.74 | -4.57 | -2.82 | -7.59 | -2025.76 |
| 34 | 41690_at | Homo sapiens mRN | 5279.81 | 1166.43 | -4.53 | -3.01 | -6.35 | -4113.38 |
| 35 | 33386_at | H1 histone family, r | 6216.48 | 1381.24 | -4.50 | -2.46 | -10.81 | -4835.24 |
| 36 | 37780_at | piccolo (presynapti | 2867.65 | 638.28 | -4.49 | -2.35 | -6.71 | -2229.37 |

# Filter genes

Filter based on expression level and on percent present calls.

# Cluster samples

# Exporting all the data from dChip

1. Use menu "Tools" $->$ "Export Expression Value".

2. Select "all genes"

3. Press "OK"

# Starting in R

Load the `affy` package.

```
> require(affy)
Loading required package: affy
Loading required package: Biobase
Loading required package: tools
Welcome to Bioconductor
 Vignettes contain introductory material.  To view,
  simply type: openVignette()
  For details on reading vignettes, see
  the openVignette help page.
Loading required package: reposTools
[1] TRUE
```

# Load the Sample Information file

```
> # remember where the data lives
> home <- 'g:/ShortCourse'
> # use the same sample info file we made for dChip
> si <- read.table(file.path(home, 'InfoFiles',
+                           'krc-sample-info.xls'),
+                    header=TRUE, sep='\t')
> si[1:5,]
           Scan.name type      Split
ALL01 CL2001011101AA  ALL Training
ALL02 CL2001011104AA  ALL Training
ALL03 CL2001011105AA  ALL Training
ALL04 CL2001011108AA  ALL Training
ALL05 CL2001011109AA  ALL Training
```

# Load dChip's array file

```
> arrays <- read.table(file.path(home, 'Output',
                       'affyShortCourse arrays.xls'),
+              header=TRUE, as.is=TRUE, sep='\t')
> # Fix the column names!
> dimnames(arrays)[[2]] <-
+     c(dimnames(arrays)[[2]][2:8], 'x')
> # Use the sample name as the row name
> dimnames(arrays)[[1]] <- arrays$Array
> # Only keep useful columns
> arrays <- arrays[, 3:6]
> # Give them sensible names
> dimnames(arrays)[[2]] <- c('MedianIntensity',
+     'PercentPresent', 'ArrayOutlier',
+     'SingleOutlier')
```

# Combine sample information

```
> # Merge sample info with dChip info
> si <- merge(si, arrays, by='row.names',
+    sort=FALSE)
> # Sigh. Fix the row names yet again.
> dimnames(si)[[1]] <- si$Row.names
> # Remove redundant columns
> si <- si[, 2:8]
> rm(arrays) # cleanup
```

# Note that the order has changed!

```
> si[1:5,]
```

|       | Scan.name      | type | Split    | MedianIntensity |
|-------|----------------|------|----------|-----------------|
| ALL01 | CL2001011101AA | ALL  | Training | 1497            |
| ALL24 | CL2001011102AA | ALL  | Test     | 1196            |
| ALL02 | CL2001011104AA | ALL  | Training | 1778            |
| ALL03 | CL2001011105AA | ALL  | Training | 1097            |
| ALL04 | CL2001011108AA | ALL  | Training | 1489            |

|       | PercentPresent | ArrayOutlier | SingleOutlier |
|-------|----------------|--------------|---------------|
| ALL01 | 48.2           | 1.648        | 0.099         |
| ALL24 | 38.3           | 3.778        | 0.432         |
| ALL02 | 49.5           | 1.450        | 0.144         |
| ALL03 | 36.8           | 3.152        | 0.375         |
| ALL04 | 38.7           | 5.299        | 0.216         |

# Specialized factors

```
> # make a factor to compare ALL vs MLL
> temp <- si$type
> temp[temp=='AML'] <- NA
> si$ALLvMLL <- factor(temp)
>
> temp <- si$type
> temp[temp=='MLL'] <- NA
> si$ALLvAML <- factor(temp)
>
> temp <- si$type
> temp[temp=='ALL'] <- NA
> si$MLLvAML <- factor(temp)
```

```
> temp <- si$type
> temp[temp=='AML'] <- 'Other'
> temp[temp=='MLL'] <- 'Other'
> si$ALLvOther <- factor(temp)
>
> temp <- si$type
> temp[temp=='ALL'] <- 'Other'
> temp[temp=='MLL'] <- 'Other'
> si$AMLvOther <- factor(temp)
>
> temp <- si$type
> temp[temp=='AML'] <- 'Other'
> temp[temp=='ALL'] <- 'Other'
> si$MLLvOther <- factor(temp)
>
> si$type <- factor(si$type)
```

```
> summary(si)
   Scan.name               type          Split
 Length:72              ALL:24      Length:72
  Class :character      AML:28       Class :character
  Mode  :character      MLL:20       Mode  :character


MedianIntensity   PercentPresent      ArrayOutlier
Min.    : 804     Min.    :28.30      Min.    : 0.253
1st Qu.:1222      1st Qu.:36.80       1st Qu.: 0.729
Median :1442      Median :41.10       Median : 1.085
Mean    :1483     Mean    :40.46      Mean    : 1.724
3rd Qu.:1727      3rd Qu.:44.83       3rd Qu.: 1.697
Max.    :3097     Max.    :49.80      Max.    :14.337
```

```
SingleOutlier        ALLvMLL        ALLvAML        MLLvAML
Min.    :0.0440      ALL :24        ALL :24        AML :28
1st Qu.:0.1610       MLL :20        AML :28        MLL :20
Median :0.2405       NA's:28        NA's:20        NA's:24
Mean    :0.2741
3rd Qu.:0.3460
Max.    :0.9520


ALLvOther      AMLvOther      MLLvOther
ALL  :24       AML  :28       MLL  :20
Other:48       Other:44       Other:52
```

# Create the phenoData object

```
> pd <- new('phenoData', pData=si, varLabels=list(
+     Scan.name='CEL file name',
+     type='Histological classification',
+     Split='Used as training or test,
+   MedianIntensity='Unnormalized median brighness,
+   PercentPresent='Percentage of present calls',
+     ArrayOutlier='Percentage of Array Outliers',
+    SingleOutlier='Percentage of Single Outliers',
+     ALLvMLL='binary classifier',
+     ALLvAML='binary classifier',
+     MLLvAML='binary classifier',
+     ALLvOtherL='binary classifier',
+     AMLvOther='binary classifier',
+     MLLvOther='binary classifier'))
```

# Create the MIAME object

MIAME = minimum information about a microarray experiment

Some of the BioConductor routines require a MIAME object, even thought they will let you submit a character string as a description.

```
> miame <- new('MIAME',
+               name='SA Armstrong',
+               lab='Lander-Golub',
+               title='MLL translocations')
```

# Read in the data from dChip

```
> temp <- read.table(file.path(home, 'Output',
+             'affyShortCourse expression.xls'),
+          header=TRUE, as.is=TRUE, sep='\t',
+          quote='', comment.char='')
> # expression data in the later columns
> data <- as.matrix(temp[, 6:77])
> # gene identifiers in the first five columns
> gi <- temp[, 1:5]
> # Use probe sets as row names
> dimnames(gi)[[1]]   <- gi$probe.set
> dimnames(data)[[1]] <- gi$probe.set
```

# Check that the order agrees

We noticed that the order of entries in the sample info file had changed when we merged it with the dChip array information. Just to be on the safe side, we should make sure that the order of the data columns matches the sample infor rows.

```
> sum(dimnames(si)[[1]] != dimnames(data)[[2]])
[1] 0
> sum(dimnames(si)[[1]] == dimnames(data)[[2]])
[1] 72
```

# Turn the dChip data into an exprSet

We can bring the dChip quantifications directly into R and turn them into an `exprSet`. Note that this avoids the memory problems by not brininging in the individual CEL files and not producing an AffyBatch.

```
> dchip <- new('exprSet',
+               exprs=data,
+               phenoData=pd,
+               annotation='hgu95av2',
+               description=miame,
+               notes='processed by KRC in dChip)
> rm(temp, data, si) # cleanup
```

# Just RMA

In order to process the data using RMA in BioConductor, we will use the `just.rma` function. This method avoids the memory problems associated with reading all the CEL files into R and keeping them around during processing. Instead, the processing is handed off to a C module that produces an `exprSet` but skips the production of an `AffyBatch` object.

# Locating CEL files in multiple directories

The default behavior for the BioConductor routines is to read all the CEL files in the current working directory. If you want to combine files from more than one location (or if you only want to use a subset of the CEL files), then you must first prepare a list of character strings that give the complete names and locations of the files you want.

```
> # CEL file location is a function of type
> cel.location <- list(ALL='CELFiles',
+        MLL='CELFiles', AML='AMLCELFiles')
> # Use the type of each file to find its location
> celdir <- cel.location[as.character(
+      pd@pData$type)]
```

```
> # Paste the '.cel' extension at the end
> celname <- paste(pd@pData$Scan.name,
+    'cel', sep='.')
> # make complete file paths for each file
> all.cel.files <- file.path(home, celdir, celname)
> # peek at the results
> all.cel.files[1:3]
[1] "g:/ShortCourse/CELFiles/CL2001011101AA.cel"
[2] "g:/ShortCourse/CELFiles/CL2001011102AA.cel"
[3] "g:/ShortCourse/CELFiles/CL2001011104AA.cel"
```

# Running Just RMA

```
> rmaData <- just.rma(filenames=all.cel.files,
+       phenoData=pd, description=miame)
[1] "Attempting to download hgu95acdf from
   http://www.bioconductor.org/packages/data/
   annotation/stable/bin/windows/contrib/2.1"
[1] "Download complete."
[1] "Installing hgu95acdf"
[1] "Installation complete"
Background correcting
Normalizing
Calculating Expression
```

# Unexpected glitches

The two quantified sets have different numbers of genes:

```
> rmaData
Expression Set (exprSet) with
12626 genes
72 samples
 phenoData object with 13 variables and 72 cases

> dchip
Expression Set (exprSet) with
12625 genes
72 samples
 phenoData object with 13 variables and 72 cases
```

# Figuring out which probe sets differ

```
> rma.ps <- dimnames(rmaData@exprs)[[1]]
> dchip.ps <- dimnames(dchip@exprs)[[1]]
> setdiff(rma.ps, dchip.ps)
 [1] "119_at"    "1215_at"    "1216_at"    "124_i_at"
 [5] "125_r_at"  "127_at"     "1301_s_at"  "1302_s_at"
 [9] "132_at"    "1429_at"    "1502_s_at"  "1829_at"
[13] "1864_at"   "1889_s_at"  "1982_s_at"  "36969_at"
[17] "383_at"    "397_at"     "412_s_at"   "426_at"
[21] "439_at"    "787_at"     "788_s_at"   "972_s_at"
[25] "985_s_at"  "997_at"
```
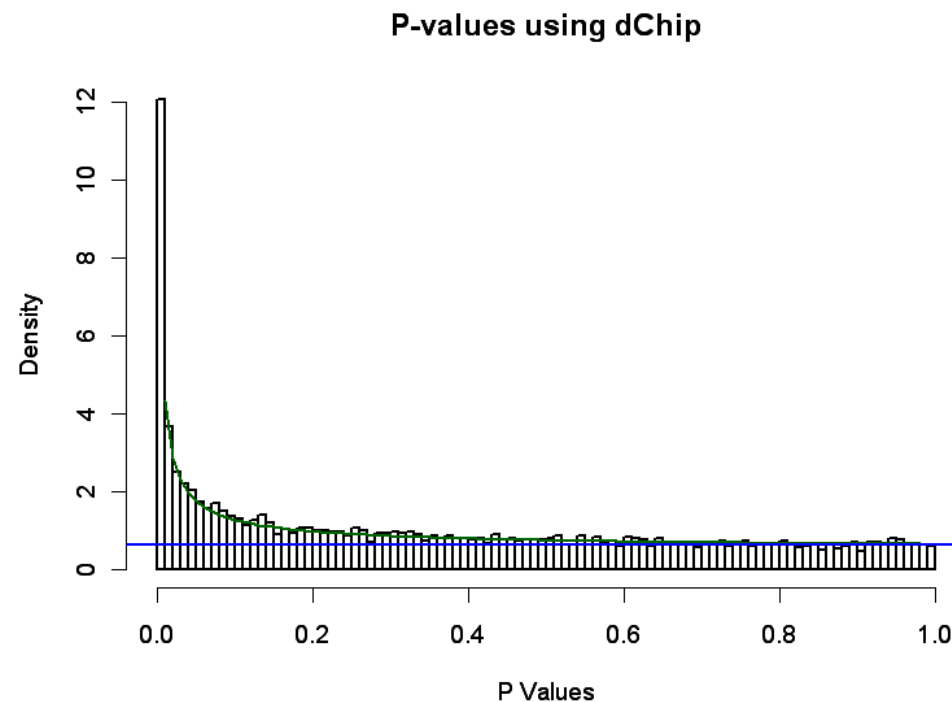
# Selecting the common probes

```
> rmaData <- rmaData[is.element(rma.ps, dchip.ps),]
> rmaData
Expression Set (exprSet) with
12600 genes
72 samples
 phenoData object with 13 variables and 72 cases
> rma.names <- dimnames(rmaData@exprs)[[1]]
> dchip@exprs <- dchip@exprs[rma.names,]
> dchip@exprs <- log(dchip@exprs, 2)
> dchip
Expression Set (exprSet) with
12600 genes
72 samples
 phenoData object with 13 variables and 72 cases
```

# Loading the ClassComparison package

```
> require(ClassComparison)
Loading required package: ClassComparison
Loading required package: splines
Loading required package: oompaBase
Loading required package: PreProcess
Creating a new generic function for 'plot' in 'PrePro
Creating a new generic function for 'print' in 'PrePi
Creating a new generic function for 'as.data.frame'
[1] TRUE
```
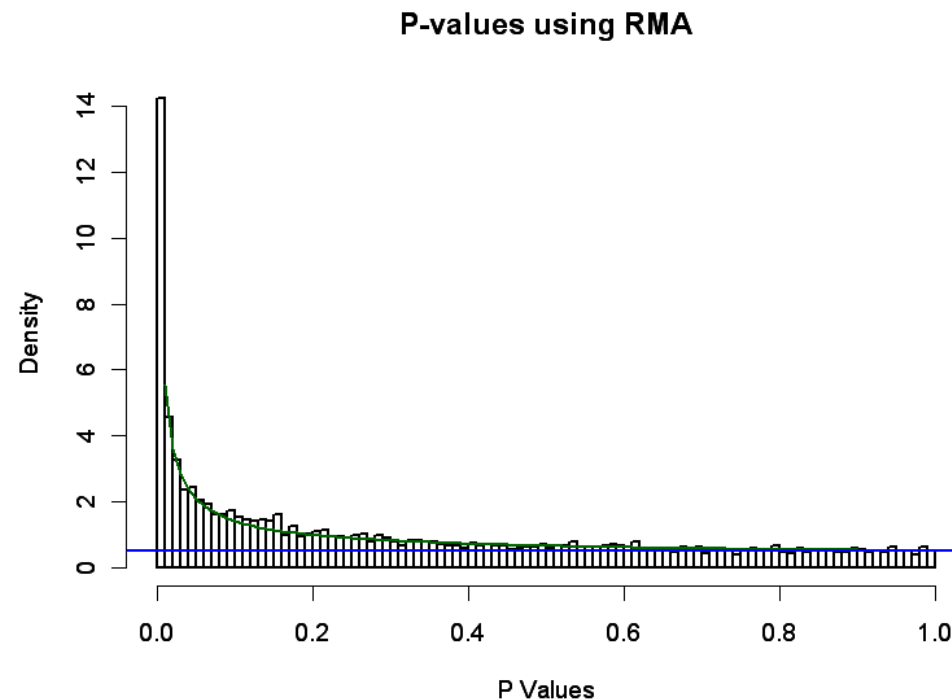
# T-test, take one

```
> notAML <- pd@pData$type != 'AML'
> dchip.t <- MultiTtest(dchip[, notAML], 'ALLvMLL')
> dchip.b <- Bum(dchip.t@p.values)
> hist(dchip.b, main='P-values using dChip')
```



P-values using dChip

# T-test, take two

```
> rma.t <- MultiTtest(rmaData[, notAML], 'ALLvMLL')
> rma.b <- Bum(rma.t@p.values)
> hist(rma.b, main='P-values using RMA')
```



P-values using RMA

# Do the two methods agree?

```
> plot(rma.t@t.statistics, dchip.t@t.statistics)
> abline(0,1, col='blue')
```