

# Ringo - R Investigation of NimbleGen Oligoarrays

Joern Toedling

## 1 Introduction

The package *Ringo* deals with the analysis of two-color oligonucleotide microarrays used in ChIP-chip projects. The package was started to facilitate the analysis of two-color microarrays from the company NimbleGen<sup>1</sup>, but the package has a modular design, such that the platform-specific functionality is encapsulated and analogous two-color tiling array platforms can also be processed. The package employs functions from other packages of the Bioconductor project (?) and provides additional ChIP-chip-specific and NimbleGen-specific functionalities.

```
> library("Ringo")
```

If you use Ringo for analyzing your data, please cite:

- Joern Toedling, Oleg Sklyar, Tammo Krueger, Jenny J Fischer, Silke Sperling, Wolfgang Huber (2007). Ringo - an R/Bioconductor package for analyzing ChIP-chip readouts. *BMC Bioinformatics*, 8:221.

## Getting help

If possible, please send questions about *Ringo* to the Bioconductor mailing list.

See <http://www.bioconductor.org/docs/mailList.html>

Their archive of questions and responses may prove helpful, too.

## 2 Reading in the raw data

For each microarray, the scanning output consists of two files, one holding the Cy3 intensities, the other one the Cy5 intensities. These files are tab-delimited text files.

The package comes with (shortened) example scanner output files, in NimbleGen's *pair* format. These files are excerpts of the ChIP-chip demo data that NimbleGen provide at their FTP site for free download. Their biological context, identification of DNA binding sites of complexes containing Suz12 in human cells, has been described before (?).

---

<sup>1</sup>for NimbleGen one-color microarrays, we recommend the Bioconductor package *oligo*

```
> exDir <- system.file("exData", package = "Ringo")
> list.files(exDir, pattern = "pair.txt")

[1] "MOD_20551_PMT1_pair.txt" "MOD_20742_PMT1_pair.txt"

> head(read.delim(file.path(exDir, "MOD_20551_PMT1_pair.txt"),
+   skip = 1))[, c(1, 4:7, 9)]
```

	IMAGE_ID	PROBE_ID	POSITION	X	Y	PM
1	20551_PMT1	SUZ100P0000021781	1	269	78	1149.33
2	20551_PMT1	SUZ100P0000021783	16	682	779	1192.00
3	20551_PMT1	SUZ100P0000021785	31	92	405	685.56
4	20551_PMT1	SUZ100P0000021787	46	219	608	562.67
5	20551_PMT1	SUZ100P0000021789	61	217	418	584.56
6	20551_PMT1	SUZ100P0000021791	76	147	406	636.22

In addition, there is a text file that holds details on the samples, including which two *pair* files belong to which sample<sup>2</sup>.

```
> read.delim(file.path(exDir, "example_targets.txt"), header = TRUE)

SlideNumber      FileNameCy3      FileNameCy5
1      Suz12 MOD_20551_PMT1_pair.txt MOD_20742_PMT1_pair.txt
      Species  Cy3  Cy5
1 Homo sapiens (human) total Suz12
```

The columns `FileNameCy3` and `FileNameCy5` hold which of the raw data files belong to which sample. The immuno-precipitated extract was colored with the Cy5 dye in the experiment, so the column `Cy5` essentially holds which antibody has been used for the immuno-precipitation, in this case one against the protein `Suz12`.

Furthermore, there is a file describing the reporter categories on the array (you might know these Spot Types files from *limma* (?)).

```
> read.delim(file.path(exDir, "spottypes.txt"), header = TRUE)

SpotType GENE_EXPR_OPTION PROBE_ID Color
1  Probe      FORWARD*      * black
2  Probe      REVERSE*      * black
3  Probe      BLOCK*        * black
4 Negative  NGS_CONTROLS*      * yellow
5  Empty      EMPTY          * white
6  H_Code     H_CODE              * red
7  V_Code     V_CODE              * blue
8  Random     RANDOM              * green
```

---

<sup>2</sup>You may have to construct such a targets file for your own data. The `scripts` directory of this package contains a script `convertSampleKeyTxt.R` as an inspiration how the file `SampleKey.txt` provided by NimbleGen could be used for this.

Reading all these files, we can read in the raw reporter intensities and obtain an object of class *RGList*, a class defined in package *limma*.

```
> exRG <- readNimblegen("example_targets.txt", "spottypes.txt",
+   path = exDir)
```

This object is essentially a list and contains the raw intensities of the two hybridizations for the red and green channel plus information on the reporters on the array and on the analyzed samples.

```
> head(exRG$R)
```

```
      MOD_20742_PMT1_pair
[1,]          613.22
[2,]          841.67
[3,]          659.56
[4,]          494.44
[5,]          469.33
[6,]          544.11
```

```
> head(exRG$G)
```

```
      MOD_20551_PMT1_pair
[1,]          1149.33
[2,]          1192.00
[3,]           685.56
[4,]           562.67
[5,]           584.56
[6,]           636.22
```

```
> head(exRG$genes)
```

	GENE_EXPR_OPTION	PROBE_ID	POSITION	X	Y	Status	ID
1	FORWARD1	SUZ100P0000021781	1	269	78	Probe	SUZ100P0000021781
2	FORWARD1	SUZ100P0000021783	16	682	779	Probe	SUZ100P0000021783
3	FORWARD1	SUZ100P0000021785	31	92	405	Probe	SUZ100P0000021785
4	FORWARD1	SUZ100P0000021787	46	219	608	Probe	SUZ100P0000021787
5	FORWARD1	SUZ100P0000021789	61	217	418	Probe	SUZ100P0000021789
6	FORWARD1	SUZ100P0000021791	76	147	406	Probe	SUZ100P0000021791

```
> exRG$targets
```

	SlideNumber	FileNameCy3	FileNameCy5
1	Suz12	MOD_20551_PMT1_pair.txt	MOD_20742_PMT1_pair.txt
		Species	Cy3 Cy5
1	Homo sapiens (human)	total	Suz12

Users can alternatively supply raw two-color ChIP-chip readouts from other platforms in *RGList* format and consecutively use *Ringo* to analyze that data.

### 3 Mapping reporters to genomic coordinates

By *reporters*, we mean the oligo-nucleotides or PCR products that have been fixated on the array for measuring the abundance of corresponding genomic fragments in the ChIP-chip experiment.

Each reporter has a unique identifier and (ideally) a unique sequence, but can, and probably does, appear in multiple copies as *features* on the array surface.

A mapping of reporters to genomic coordinates is usually provided by the array manufacturer, such as in NimbleGen's \*.POS files. If the reporter sequences are provided as well, you may consider to perform a custom mapping of these sequences to the genome of interest, using alignment tools such as *Exonerate* (?) or functions provided by the Bioconductor package *Biostrings* (?).

Such a re-mapping of reporters to the genome can sometimes be necessary, for example when the array has been designed on an outdated assembly of the genome. Re-mapping also provides the advantage that you can allow non-perfect matches of reporters to the genome, if desired.

Once reporters have been mapped to the genome, this mapping needs to be made available to the data analysis functions. While a *data.frame* may be an obvious way of representing such a mapping, repeatedly extracting sub-sets of the data frame related to a genomic region of interest turns out to be too slow for practical purposes. *Ringo*, similar to the Bioconductor package *tilingArray*, employs an object of class *probeAnno* to store the mapping between reporters on the microarray and genomic positions. Per chromosome, the object holds four vectors of equal length and ordering that specify at which genomic positions reporter matches start and end, what identifiers or indices these reporters have in the intensities data, and whether these reporters match uniquely to the genomic positions.

```
> load(file.path(exDir, "exampleProbeAnno.rda"))
> ls(exProbeAnno)
```

```
[1] "9.end"      "9.index"    "9.start"    "9.unique"
```

```
> show(exProbeAnno)
```

A 'probeAnno' object holding the mapping between reporters and genomic positions.

Chromosomes: 9

Microarray platform: NimbleGen MOD\_2003-12-05\_SUZ12\_1in2

Genome: H.sapiens (hg18)

```
> head(exProbeAnno["9.start"])
```

```
[1] 531857 861532 1269645 2497214 2685042 3795381
```

```
> head(exProbeAnno["9.end"])
```

```
[1] 531916 861591 1269704 2497273 2685101 3795440
```

The function `posToProbeAnno` allows generation of a valid *probeAnno* object, either from a file that corresponds to a NimbleGen POS file or from a *data.frame* objects that holds the same information. The package's `scripts` directory contains the script `mapReporter-sWithBiostrings.R`, which shows how to use *Biostrings* for mapping the reporter sequences of the provided example data. and some Perl scripts that allow the conversion of multiple output files from common alignment tools such as Exonerate into one file that corresponds to a POS file. The function `validObject` can be used to perform a quick check whether a generated *probeAnno* object will probably work with other *Ringo* functions.

## 4 Quality assessment

The `image` function allows us to look at the spatial distribution of the intensities on a chip. This can be useful to detect obvious artifacts on the array, such as scratches, bright spots, finger prints etc. that might render parts or all of the readouts useless.

```
> par(mar = c(0.01, 0.01, 0.01, 0.01), bg = "black")
> image(exRG, 1, channel = "green", mycols = c("black", "green4",
+       "springgreen"))
```

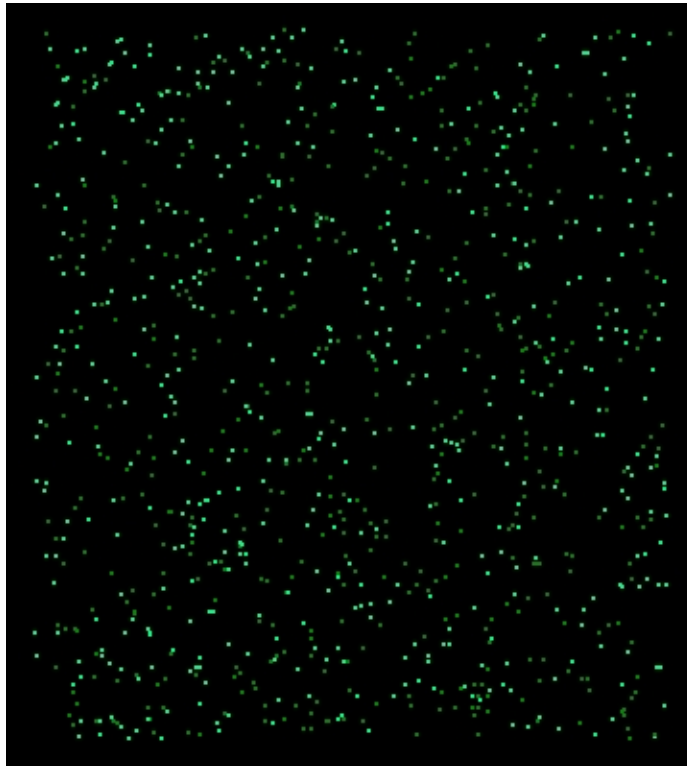


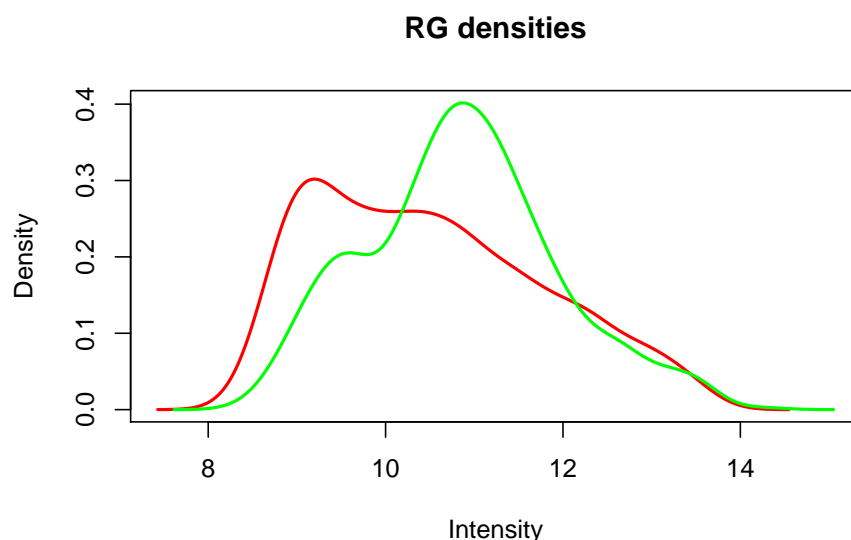
Figure 1: *Spatial distribution of raw reporter intensities laid out by the reporter position on the microarray surface.*

See figure 1 for the image. Since the provided example data set only holds the intensities for reporters mapped to the forward strand of chromosome 9, the image only shows the few

green dots of these reporters' positions. We see, however, that these chromosome 9 reporters are well distributed over the whole array surface rather than being clustered together in one part of the array.

It may also be useful to look at the absolute distribution of the single-channel densities. *limma*'s function `plotDensities` may be useful for this purpose.

```
> plotDensities(exRG)
```



In addition, the data file loaded above also contains a *GFF* (*General Feature Format*) file of all transcripts on human chromosome 9 annotated in the [Ensembl](#) database (release 46, August 2007). The script `retrieveGenomicFeatureAnnotation.R` in the package's scripts directory contains example source code showing how the Bioconductor package *biomaRt* can be used to generate such an annotated genome features *data.frame*.

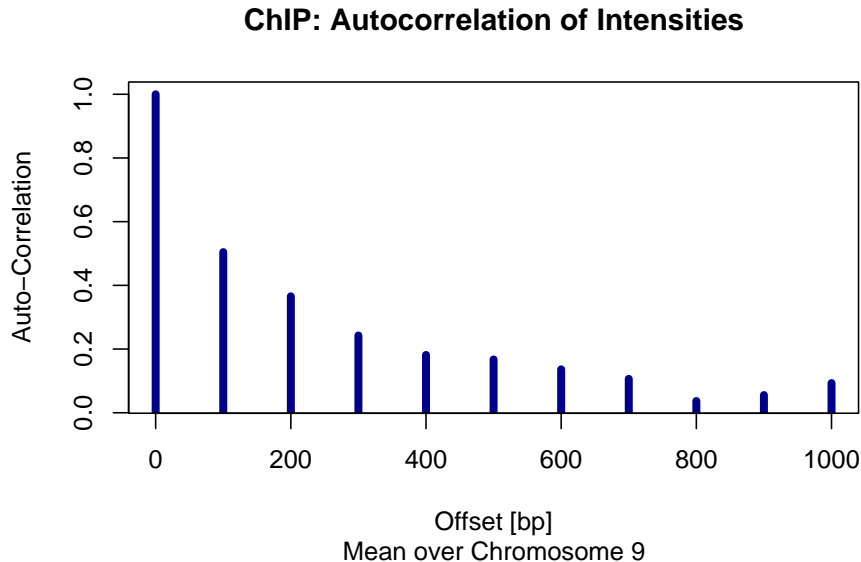
```
> head(exGFF[, c("name", "symbol", "chr", "strand", "start", "end")])
```

	name	symbol	chr	strand	start	end
1	ENST00000382507		9	1	1056	1620
2	ENST00000326592		9	-1	5841	19587
3	ENST00000382502		9	-1	24394	25841
4	ENST00000355822		9	-1	24921	25277
5	ENST00000382501		9	-1	25028	25856
6	ENST00000382500	FOXD4	9	-1	105991	108951

To assess the impact of the small distance between reporters on the data, one can look at the autocorrelation plot. For each base-pair lag  $d$ , it is assessed how strong the intensities of reporters at genomic positions  $x + d$  are correlated with the probe intensities at positions  $x$ .

The computed correlation is plotted against the lag  $d$ .

```
> exAc <- autocor(exRG, probeAnno = exProbeAnno, chrom = "9", lag.max = 1000)
> plot(exAc)
```



We see some auto-correlation between probe position up to 800 base pairs apart. Since the sonicated fragments that are hybridized to the array have an average size in the range of up to 1000 bp, such a degree of auto-correlation up to this distance can be expected.

## 5 Preprocessing

Following quality assessment of the raw data, we perform normalization of the probe intensities and derive fold changes of reporters' intensities in the enriched sample divided by their intensities in the non-enriched *input* sample and take the (generalized) logarithm of these ratios.

We use the variance-stabilizing normalization (?) or probe intensities and generate an `ExpressionSet` object of the normalized probe levels.

```
> exampleX <- preprocess(exRG)
> sampleNames(exampleX) <- with(exRG$targets, paste(Cy5, "vs",
+   Cy3, sep = "_"))
> print(exampleX)
```

```
ExpressionSet (storageMode: lockedEnvironment)
assayData: 991 features, 1 samples
  element names: exprs
phenoData
  sampleNames: Suz12_vs_total
  varLabels and varMetadata description:
    SlideNumber: NA
    FileNameCy3: NA
```

```

....: ...
Cy5: NA
(6 total)
additional varMetadata: varLabel
featureData
  featureNames: SUZ100P0000021781, SUZ100P0000021783, ..., SUZ100P0000058079 (991 total)
  fvarLabels and fvarMetadata description: none
experimentData: use 'experimentData(object)'
Annotation:

```

Among the provided alternative preprocessing options is also the Tukey-biweight scaling procedure that NimbleGen have used to scale ChIP-chip readouts so that the data is centered on zero.

```

> exampleX.NG <- preprocess(exRG, method = "nimblegen")
> sampleNames(exampleX.NG) <- sampleNames(exampleX)

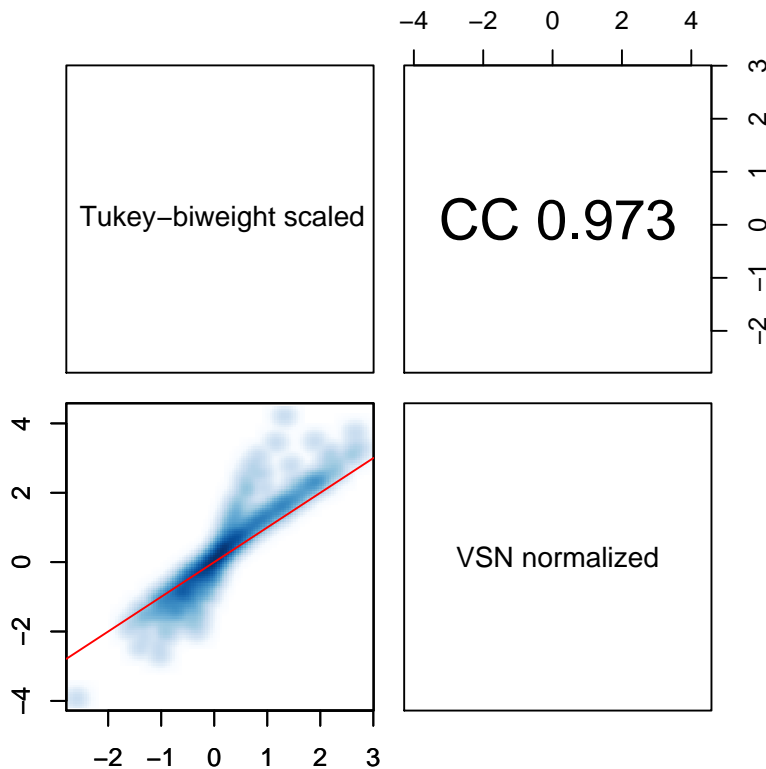
```

The effects of different preprocessing procedures on the data, can be assessed using the `corPlot` function.

```

> corPlot(cbind(exprs(exampleX), exprs(exampleX.NG)), grouping = c("VSN normalized",
+   "Tukey-biweight scaled"))

```





The same function can also be used to assess the correlation between biological and technical replicates among the microarray samples.

## 6 Visualize intensities along the chromosome

```
> chipAlongChrom(exampleX, chrom = "9", xlim = c(34318000, 34321000),
+   ylim = c(-2, 4), probeAnno = exProbeAnno, gff = exGFF)
```

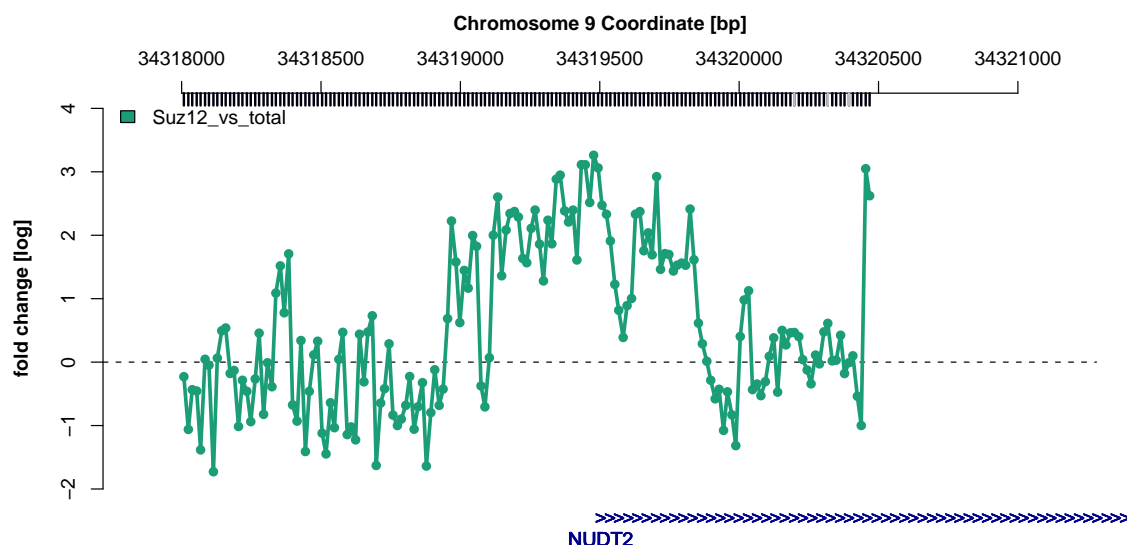


Figure 2: *Normalized probe intensities around the TSS of the *Nudt2* gene.*

See the result in figure 2.

## 7 Smoothing of probe intensities

Since the response of reporters to the same amount of hybridized genome material varies greatly, due to probe GC content, melting temperature, secondary structure etc., it is suggested to do a smoothing over individual probe intensities before looking for ChIP-enriched regions.

Here, we slide a window of 800 bp width along the chromosome and replace the intensity at a genomic position  $x_0$  by the median over the intensities of those reporters inside the window that is centered at  $x_0$ .

```
> smoothX <- computeRunningMedians(exampleX, probeAnno = exProbeAnno,
+   modColumn = "Cy5", allChr = "9", winHalfSize = 400)
> sampleNames(smoothX) <- paste(sampleNames(exampleX), "smoothed")

> chipAlongChrom(exampleX, chrom = "9", xlim = c(34318000, 34321000),
+   ylim = c(-2, 4), probeAnno = exProbeAnno, gff = exGFF)
```

```
> chipAlongChrom(smoothX, chrom = "9", xlim = c(34318000, 34321000),
+   probeAnno = exProbeAnno, itype = "l", ilwd = 4, paletteName = "Spectral",
+   add = TRUE)
```

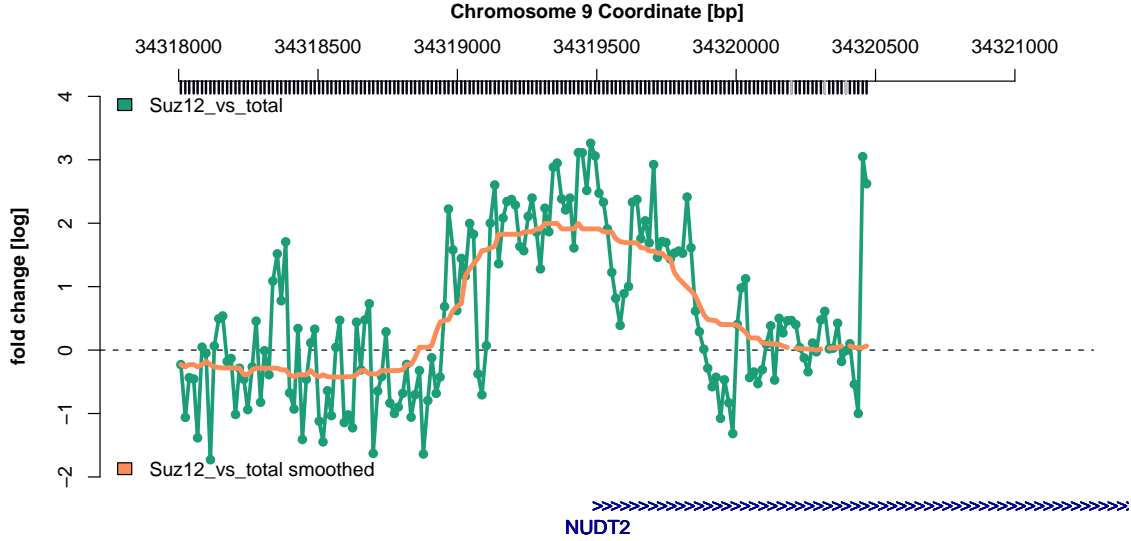


Figure 3: *Normalized and smoothed probe intensities around the TSS of the *Nudt2* gene.*

See the smoothed probe levels in figure 3.

## 8 Finding ChIP-enriched regions

To identify antibody-enriched genomic regions, we require the following:

- smoothed intensities of reporters mapped to this region are exceed a certain threshold  $y_0$
- the region contains at least three probe match positions
- each affected position is less than a defined maximum distance  $d_{max}$  apart from another affected position in the region (we require a certain probe spacing to have confidence in detected peaks<sup>3</sup>)

For setting the threshold  $y_0$ , one has to assess the expected (smoothed) probe levels in non-enriched genomic regions, i.e. the *null distribution* of probe levels. In a perfect world, we could use a log ratio of 0 as definite cut-off. In this case the “enriched” DNA and the input DNA sample would be present in equal amounts, so no antibody-bound epitope, could be found at this genomic site. In practice, there are some reasons why zero may be

<sup>3</sup>Note that the term “peak”, while commonly used in ChIP-chip context, is slightly misleading and the term “ChIP-enriched region”, or “cher” in shorthand, is more appropriate. Within such regions the actual signal could show two or more actual signal peaks or none at all (long plateau).

a too naive cut-off for calling a probe-hit genomic site *enriched* in our case. See ? for an extensive discussion on problematic issues with ChIP-chip experiments. We will just briefly mention a few issues here. For once, during the immuno-precipitation, some non-antibody-bound regions may be pulled down in the assay and consequently enriched or some enriched DNA may cross-hybridize to other reporters. Furthermore, since genomic fragments after sonication are mostly a lot larger than the genomic distance between two probe-matched genomic positions, auto-correlation between reporters certainly is existent. Importantly, different reporters measure the same DNA amount with a different efficiency even after normalizing the probe levels, due to sequence properties of the probe, varying quality of the synthesis of reporters on the array and other reasons. To ameliorate this fact, we employ the sliding-window smoothing approach.

The aforementioned issues make it difficult to come up with a reasonable estimate for the null distribution of smoothed probe levels in non-enriched genomic regions. See Figure 4 for the two histograms. We present one way (out of many) for objectively choosing the cutoff  $y_0$ . The histograms suggest the smoothed reporter levels follow a mixture of two distributions, one being the null distribution of non-affected reporters and the other one the alternative one for the smoothed reporter values in H3K4me3-ChIP enriched regions. We assume the null distribution is symmetric and its mode is the one close to zero in the histogram. By mirroring its part left of the mode over the mode, we end up with an estimated null distribution. For the alternative distribution, we only assume that it is stochastically larger than the null distribution and that its mass to the left of the estimated mode of the null distribution is negligible. We estimate an upper bound  $y_0$  for values arising from the null distribution and conclude that smoothed probe levels  $y > y_0$  are more likely to arise from the H3K4me3 enrichment distribution than from the null distribution. These estimates are indicated by red vertical lines in the histograms.

```
> (y0 <- apply(exprs(smoothX), 2, upperBoundNull))

Suz12_vs_total smoothed
0.7392157

> par(font.lab = 2, mar = c(4, 4, 1, 1))
> hist(exprs(smoothX)[, 1], n = 25, xlab = "Smoothed reporter intensities [log]",
+      xlim = c(-2, 2), main = NA, col = brewer.pal(8, "Set2")[1])
> abline(v = y0[, 1], col = "red")
```

Since antibodies vary in their efficiency to bind to their target epitope, we suggest to obtain a different threshold for each antibody. In the example data, however, we have only one antibody against Suz12.

While this threshold worked well for us, we do not claim this way to be a gold standard for determining the threshold. In particular, it does not take into account the auto-correlation between near-by reporters. See ? for a more sophisticated algorithm that does take it into account.

```
> chersX <- findChersOnSmoothed(smoothX, probeAnno = exProbeAnno,
+   thresholds = y0, allChr = "9", distCutOff = 600, cellType = "human")
> chersX <- relateChers(chersX, exGFF)
> chersXD <- as.data.frame.cherList(chersX)
```

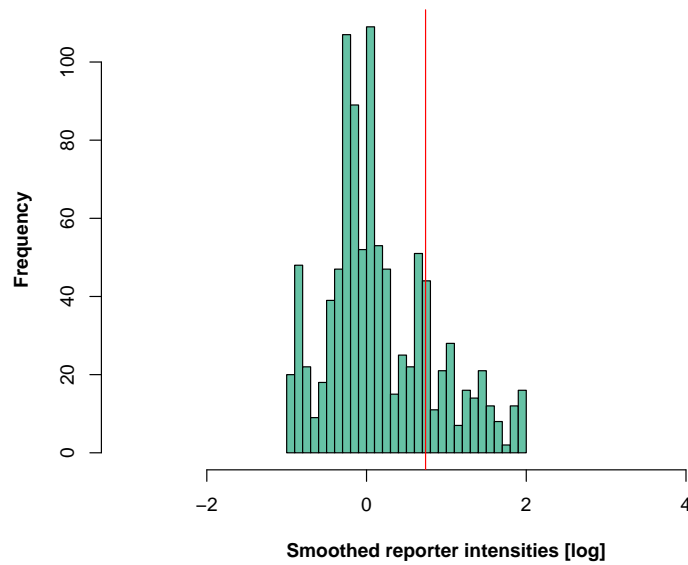


Figure 4: Histograms of reporter intensities after smoothing of reporter level. The red vertical line is the cutoff values suggested by the histogram.

```
> chersXD[order(chersXD$maxLevel, decreasing = TRUE), ]
```

	name	chr	start	end	cellType
1	human.Suz12_vs_total smoothed.chr9.cher1	9	34319028	34319854	human
3	human.Suz12_vs_total smoothed.chr9.cher3	9	34580420	34582384	human
2	human.Suz12_vs_total smoothed.chr9.cher2	9	34579444	34579760	human

	antibody
1	Suz12_vs_total smoothed
3	Suz12_vs_total smoothed
2	Suz12_vs_total smoothed

	features
1	ENST00000379158 ENST00000379154 ENST00000379155 ENST00000346365 ENST00000337747
3	ENST00000378980 ENST00000351266
2	ENST00000378980 ENST00000351266

	maxLevel	score
1	1.9958907	52.1036588
3	1.5341497	47.8557696
2	0.7882005	0.5865187

Note that in *Ringo* functions, “ChIP-enriched region” is abbreviated to “cher”.

One characteristic of enriched regions that can be used for sorting them is the element `maxLevel`, that is the highest smoothed probe level in the enriched region. Alternatively, one can sort by the `score`, that is the sum of smoothed probe levels minus the threshold. It is a discretized version of the area under the curve with the baseline being the threshold.

```
> plot(chersX[[1]], smoothX, probeAnno = exProbeAnno, gff = exGFF,
+       paletteName = "Spectral")
```

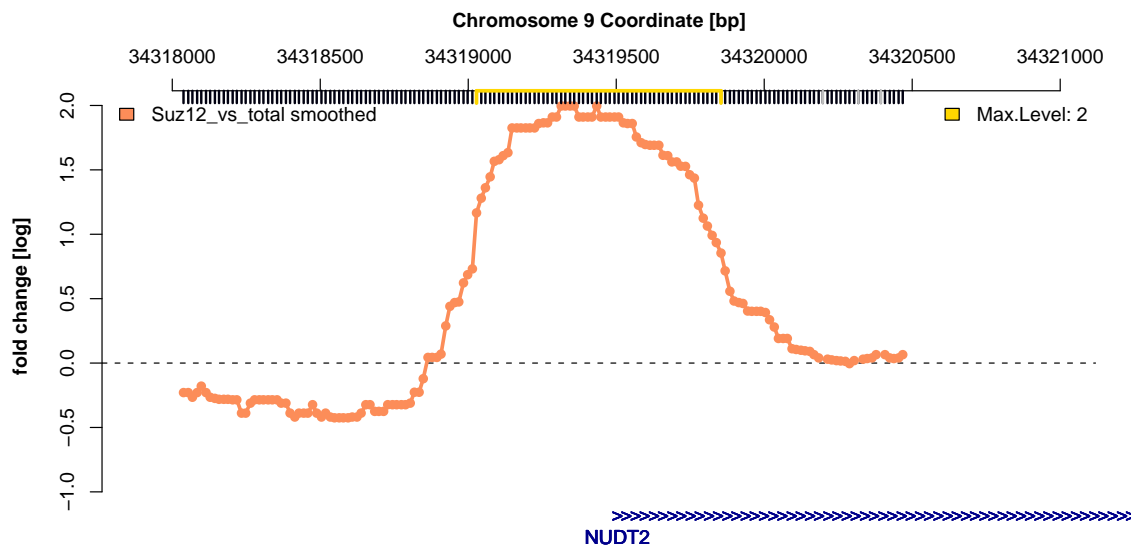


Figure 5: *One of the identified Suz12-antibody enriched regions on chromosome 9.*

Figure 5 displays an identified enriched region, which is located upstream of the *Nudt2* gene. This ChIP-enriched region was already obvious in plots of the normalized data (see Figure 3). While it is reassuring that our method recovers it as well, a number of other approaches would undoubtedly have reported it as well.

## 9 Concluding Remarks

The package *Ringo* aims to facilitate the analysis ChIP-chip readouts. We constructed it during the analysis of a ChIP-chip experiment for the genome-wide identification of modified histone sites on data gained from NimbleGen two-color microarrays. Analogous two-color microarray platforms, however, can also be processed. Key functionalities of *Ringo* are data read-in, quality assessment, preprocessing of the raw data, and visualization of the raw and preprocessed data. The package contains a heuristic algorithm for the detection of for ChIP-enriched genomic regions, too. While this algorithm worked quite well on our data, we do not claim it to be the definite algorithm for that task.

This vignette was generated using the following package versions:

- R version 2.8.0 (2008-10-20), i386-pc-mingw32
- Locale: LC\_COLLATE=English\_United States.1252;LC\_CTYPE=English\_United States.1252;LC\_MONETARY=English\_United States.1252;LC\_NUMERIC=English\_United States.1252;LC\_TIME=English\_United States.1252
- Base packages: base, datasets, graphics, grDevices, methods, splines, stats, tools, utils
- Other packages: affy 1.20.0, annotate 1.20.0, AnnotationDbi 1.4.0, Biobase 2.2.0, genefilter 1.22.0, geneplotter 1.20.0, lattice 0.17-15, limma 2.16.0, RColorBrewer 1.0-2, Ringo 1.6.0, SparseM 0.78, survival 2.34-1, vsn 3.8.0, xtable 1.5-4
- Loaded via a namespace (and not attached): affyio 1.10.0, DBI 0.2-4, grid 2.8.0, KernSmooth 2.22-22, preprocessCore 1.4.0, RSQLite 0.7-0

## Acknowledgments

I thank Wolfgang Huber Oleg Sklyar, Tammo Krüger, Richard Bourgon, and Matt Ritchie for source code contributions to and lots of helpful suggestions on Ringo, Todd Richmond and NimbleGen Systems, Inc. for providing us with the example ChIP-chip data.

This work was supported by the European Union (FP6 HeartRepair, LSHM-CT-2005-018630).

## References