# Quality Report for Affymetrix Microarray Experiment Dilution

January 10, 2007

## Contents

This is a quality assessment report for the dataset *Dilution*. The data are comprised of 4 arrays, of type `HG_U95Av2`.

For details on the software packages that were used to produce this report see Section 4.

## 1 The quality metrics recommended by Affymetrix

Affymetrix recommends a number of quality metrics that can be calculated for each array.

- Average background intensity, scale factors and percent of genes called present. These are shown in Table 1. The values should be similar across arrays. In the presented data, the ratio of the largest to the smallest value of average background is 1.737. Since this ratio is less than 3 there is unlikely to be

1

a problem. Among the scale factors, the ratio of the maximum to the minimum value is 2.066. Since this ratio is less than 3 there is unlikely to be a problem. For the percent present calls, it is 1.021. Since this ratio is less than 3 there is unlikely to be a problem.

- Ratios of hybridization efficiency between probes at the 3' and 5' ends of some control probe sets. These are displayed in Table 2. They should all be less than 3.

- External control probes. The protocols suggest that labelled cRNAs be added during sample preparation. These are BioB, BioC, BioD and CreX and are derived from Bacillus subtiliis. Nothing else should bind to their probesets. The results for these quantities are reported in Table 3. It is intended that BioB be spiked in at the lower limit of detection and that BioC, BioD and CreX be spiked in at higher concentrations. If BioB is routinely absent, then there may be a problem with sensitivity.

| | AvBg | ScaleF | PerCPres |
|---|---|---|---|
| 20A | 94.25 | 0.89 | 48.79 |
| 20B | 63.64 | 1.27 | 49.82 |
| 10A | 80.09 | 1.14 | 49.38 |
| 10B | 54.26 | 1.85 | 49.76 |

Table 1: Average background, scale factor and percent present calls.

| | a | b | c | d |
|---|---|---|---|---|
| 20A | 0.70 | 0.44 | 0.13 | −0.06 |
| 20B | 0.72 | 0.35 | 0.18 | −0.01 |
| 10A | 0.87 | 0.43 | 0.21 | 0.42 |
| 10B | 0.93 | 0.57 | 0.27 | 0.11 |

Table 2: 3'/5' ratios. a) HSAC07/X00351 3'/5' b) HUMGAPDH/M33197 3'/5' c) HSAC07/X00351 3'/M d) HUMGAPDH/M33197 3'/M.

These quality metrics are also summarized in Figure 1. Any metric that is shown in red is out of the manufacturer's specified boundaries and suggests a potential problem.
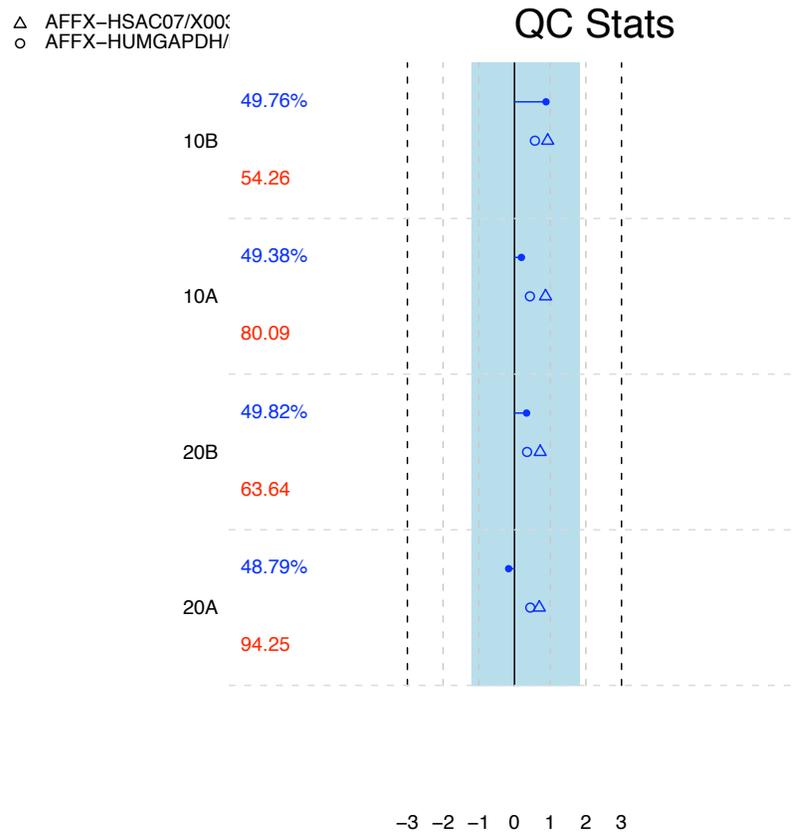
Figure 1: Quality metrics overview diagnostic plot.

|      | BioBCall | BioB   | BioC  | BioDn  | CreX  |
|------|----------|--------|-------|--------|-------|
| 20A  | P        | 11.641 | 7.98  | 11.921 | 7.88  |
| 20B  | P        | 11.231 | 7.683 | 11.729 | 7.752 |
| 10A  | P        | 11.893 | 8.192 | 12.208 | 7.888 |
| 10B  | P        | 12.183 | 8.644 | 12.506 | 8.358 |

Table 3: BioB and friends

The quality metrics reported in this Section and Figure 1 were generated using the simpleaffy package. For further information, we recommend the documentation and vignettes in the simpleaffy package.

## 2 Per array intensity distributions

### 2.1 Before normalization

The quality metrics in this section look at the distribution of the (raw, unnormalized) feature intensities for each array. Figure 2 shows density estimates (histograms), and Figure 3 presents boxplots of the same data. Arrays whose distributions are very different from the others should be considered for possible problems.

### 2.2 After normalization

$MA$-plots are useful for pairwise comparisons between arrays. $M$ and $A$ are defined as

$$
\begin{aligned}
M &= \log_2(X_1) - \log_2(X_2) = \log_2 \frac{X_1}{X_2}, \\
A &= \frac{1}{2}\left(\log_2(X_1) + \log_2(X_2)\right) = \log_2 \sqrt{X_1 X_2},
\end{aligned}
$$

where $X_1$ and $X_2$ are the vectors of normalized intensities of two arrays, on the original data scale (i. e. not logarithm-transformed).

For the $MA$-plots shown in Figure 4, the data were background corrected and normalized, but not summarized (so there is one value per probe, not one value per probeset). Rather than comparing each array to every other array, here we compare each array to a single median "pseudo"-array.

Typically, we expect the mass of the distribution in an $MA$-plot to be concentrated along the $M = 0$ axis, and there should be no trend in the mean of $M$ as a function of $A$.
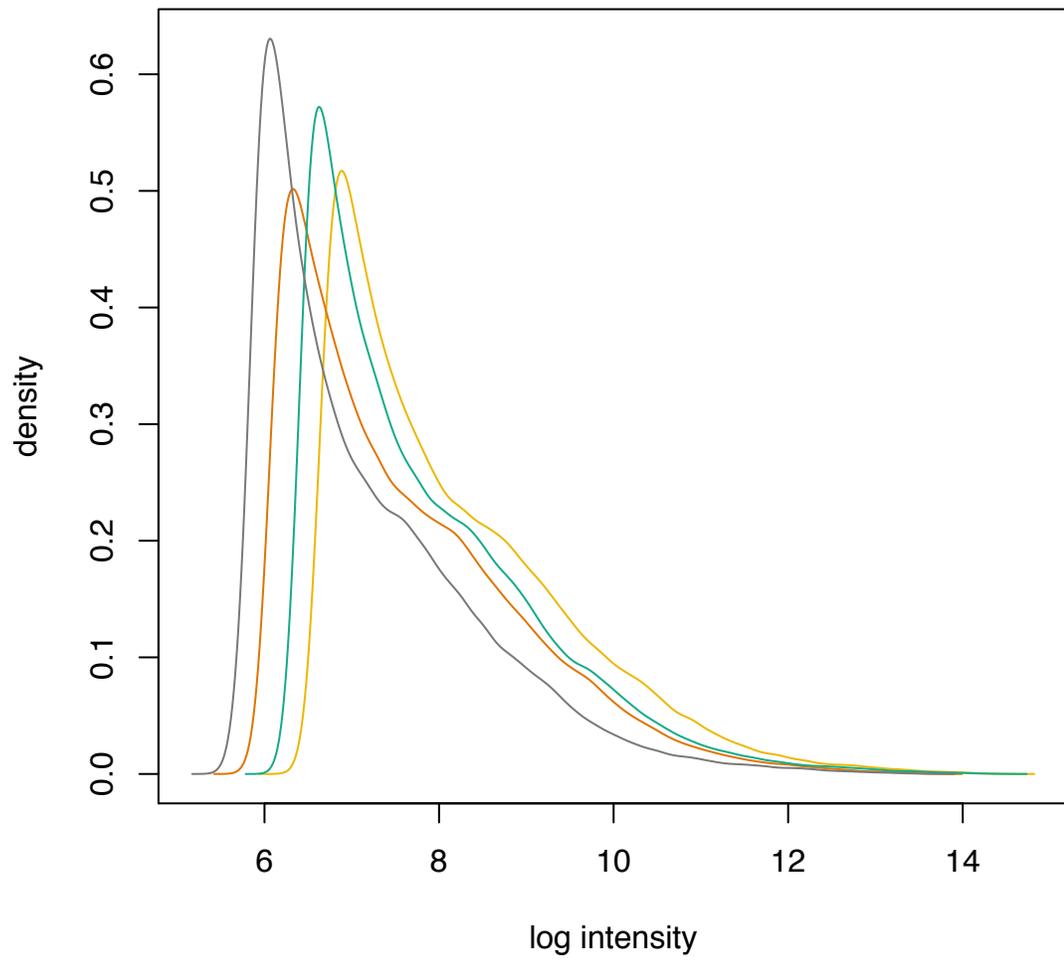
4

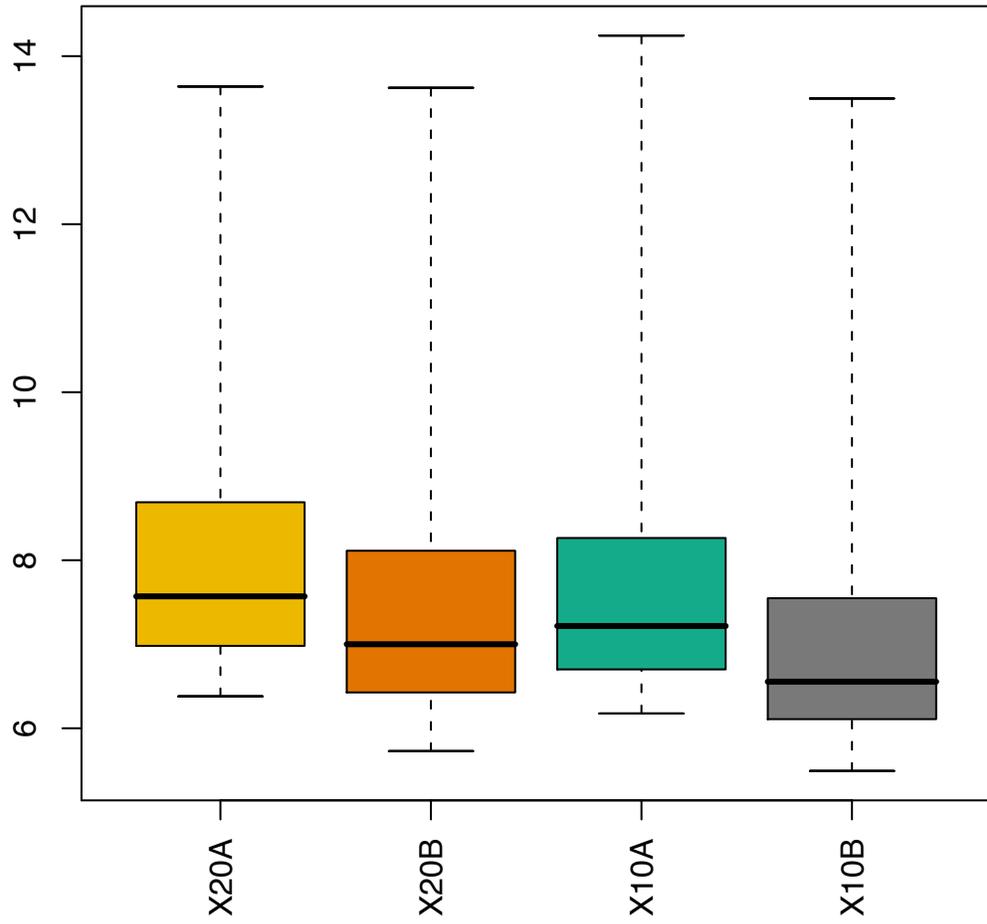Figure 2: Density estimates (histograms) for arrays 20A, 20B, 10A, 10B.

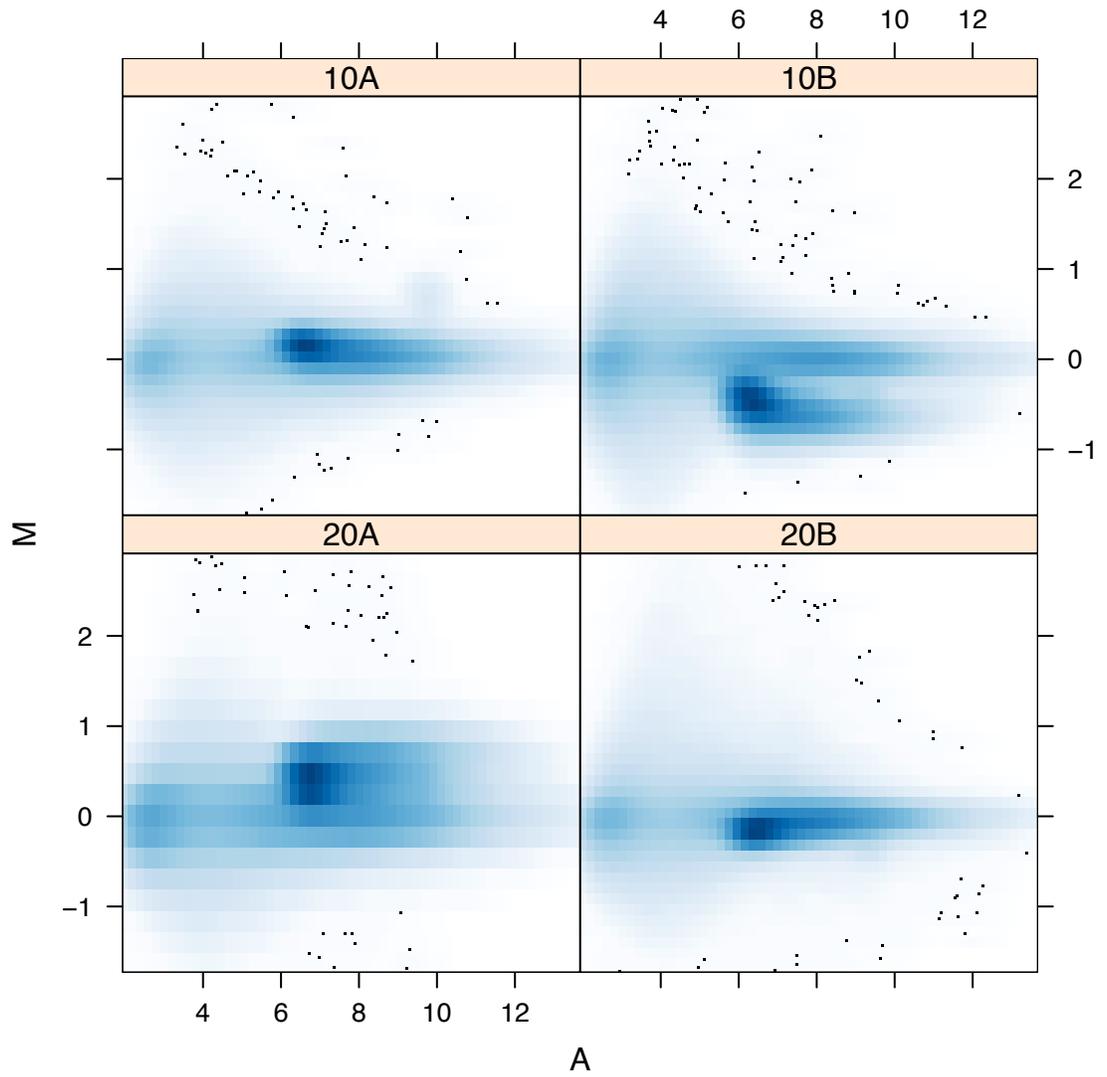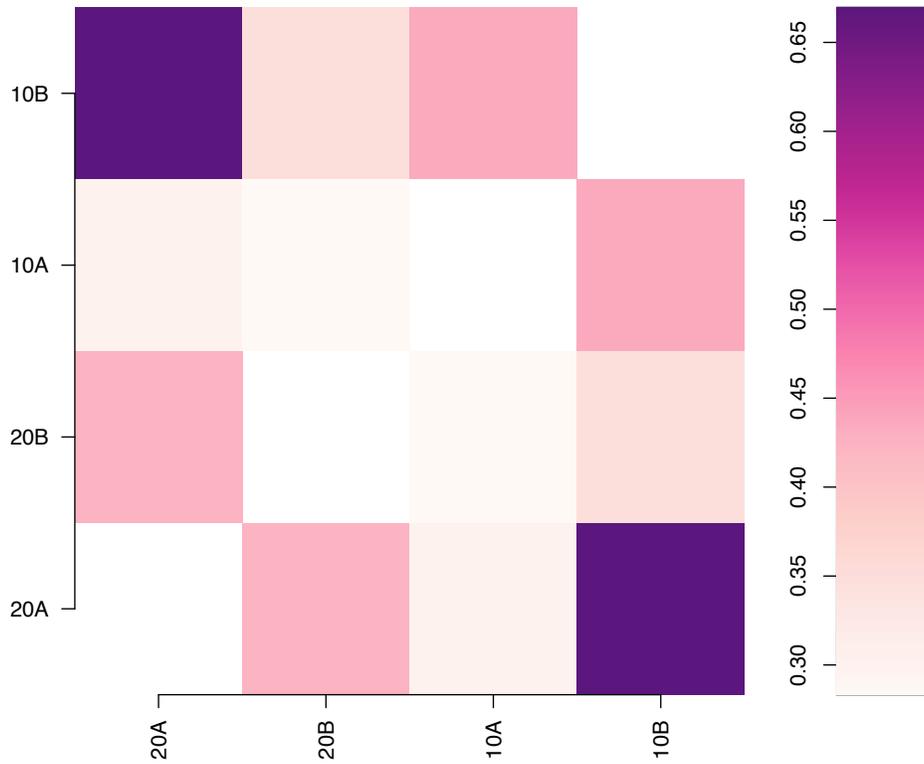Figure 3: Boxplots for arrays 20A, 20B, 10A, 10B.

Figure 4: MA plots. A *reference array* array is calculated from the median across arrays, and for each array $M$ and $A$ values are calculated for the comparison to that reference.

Figure 5: Pairwise differences between arrays, computed as the median absolute deviation (MAD) of the differences of the $M$-values.

Note that a bigger width of the plot of the $M$-distribution at the lower end of the $A$ scale does not necessarily imply that the variance of the $M$-distribution is larger at the lower end of the $A$ scale: the visual impression might simply be caused by the fact that there is more data at the lower end of the $A$ scale. To visualize whether there is a trend in the variance of $M$ as a function of $A$, consider plotting $M$ versus rank($A$).

# 3    Between array comparisons

Figure 5 shows a false color display of between arrays distances, computed as the MAD of the $M$-values of each pair of arrays.

$$d_{ij} = c \cdot \underset{m}{\mathrm{median}} \left| x_{mi} - x_{mj} \right|.$$

Here, $x_{mi}$ is the normalized intensity value of the $m$-th probe on the $i$-th array, on the original data scale. $c = 1.4826$ is a constant factor that ensures consistency with the empirical variance for Normally distributed data (see manual page of the *mad* function in R).

This plot can serve to detect outlier arrays. Consider the following decomposition of $x_{mi}$:

$$x_{mi} = z_m + \beta_{mi} + \varepsilon_{mi}, \tag{1}$$

where $z_m$ is the probe effect for probe $m$ (the same across all arrays), $\varepsilon_{mi}$ are i.i.d. random variables with mean zero and $\beta_{mi}$ is such that for any array $i$, the majority of values $\beta_{mi}$ are negligibly small (i. e. close to zero). $\beta_{mi}$ represents differential expression effects. In this model, all values $d_{ij}$ are (in expectation) the same, namely $\sqrt{2}$ times the standard deviation of $\varepsilon_{mi}$. Arrays whose distance matrix entries are way different give cause for suspicion.

## 4    Other plots (degradation and affyPLM)

In this section we present diagnostic plots based on tools provided in the affyPLM package.

In Figure 6 a RNA digestion plot is computed. In this plot each array is represented by a single line. It is important to identify any array(s) that has a slope which is very different from the others. The indication is that the RNA used for that array has potentially been handled quite differently from the other arrays.

Figure 7 is a Normalized Unscaled Standard Error (NUSE) plot. Low quality arrays are those that are significantly elevated or more spread out, relative to the other arrays. NUSE values are not comparable across data sets.

Figure 8 is a Relative Log Expression (RLE) plot and an array that has problems will either have larger spread, or will not be centered at M = 0, or both.
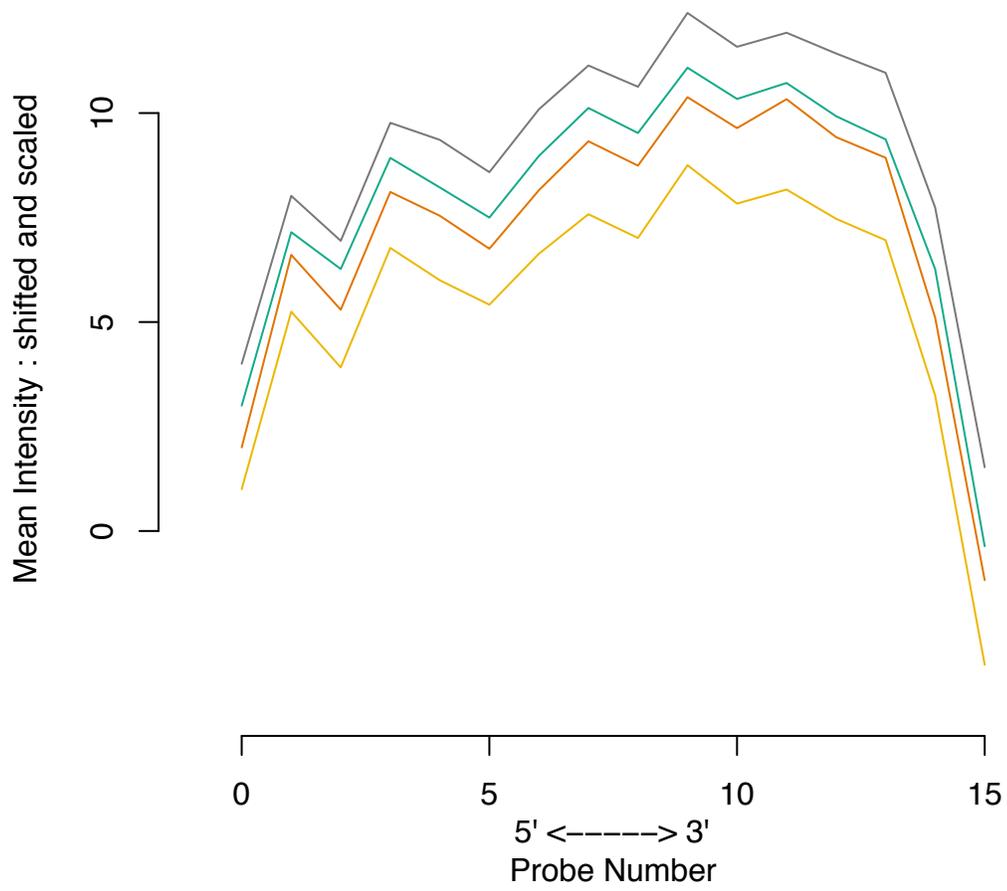
## Acknowledgements

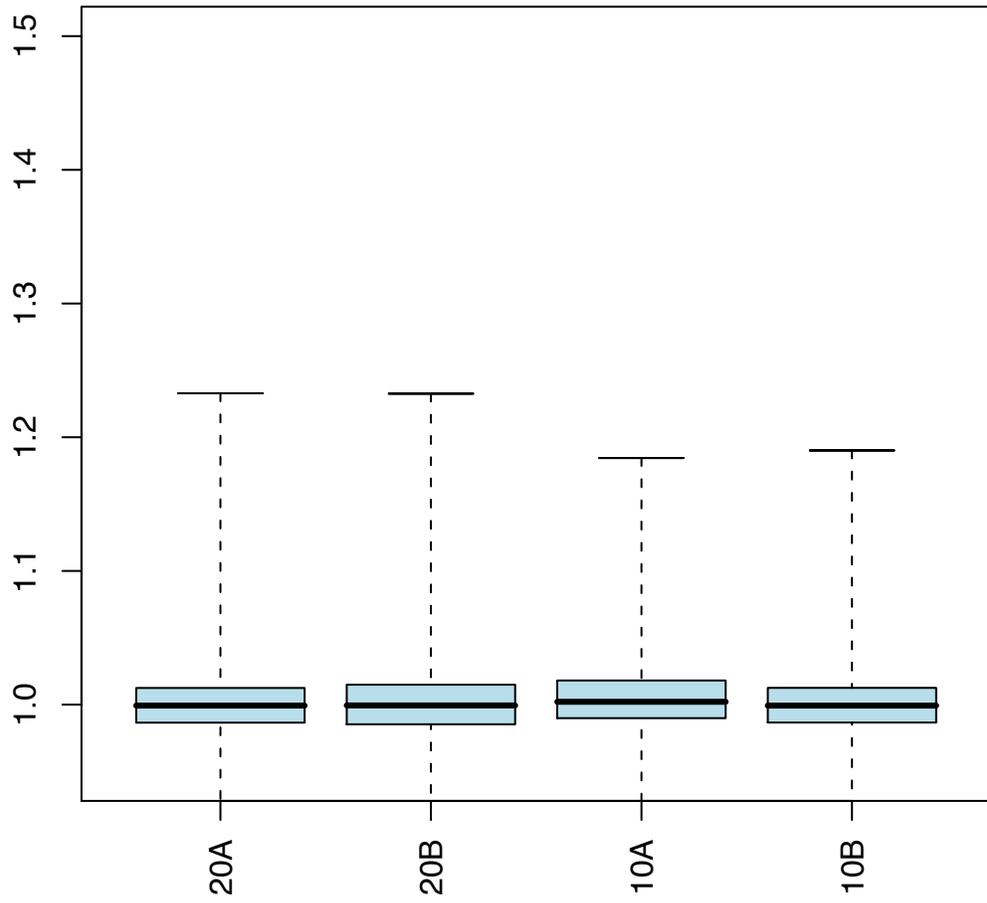Figure 6: RNA digestion / degradation plots for arrays 20A, 20B, 10A, 10B.
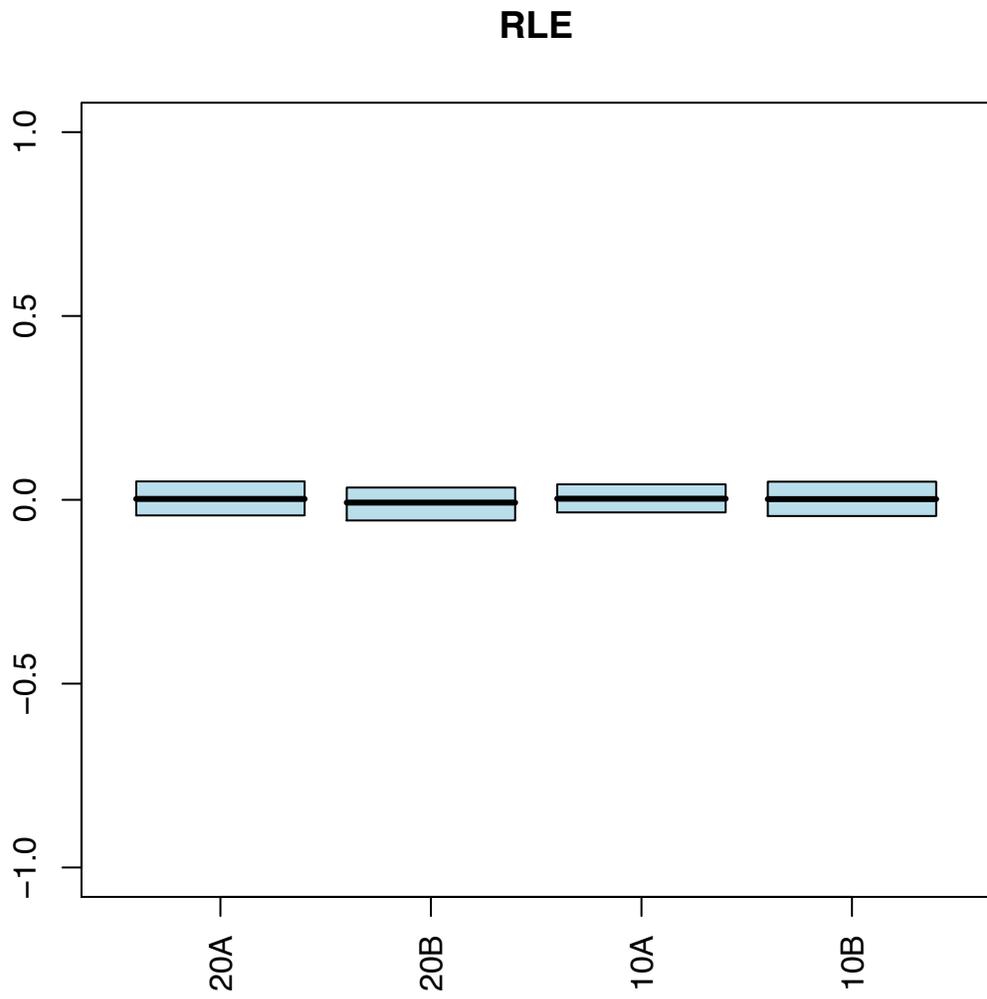
**NUSE**



Figure 7: NUSE plot.

**RLE**



Figure 8: RLE plot.

## SessionInformation:

- R version 2.5.0 Under development (unstable) (2006-11-27 r40032), `i386-apple-darwin8.8.1`

- Locale: `C`

- Base packages: base, datasets, grDevices, graphics, methods, splines, stats, tools, utils

- Other packages: Biobase 1.13.30, RColorBrewer 0.2-3, affy 1.13.12, affy-PLM 1.11.13, affyQCReport 1.13.15, affydata 1.11.1, affyio 1.3.1, annotate 1.13.3, gcrma 2.7.1, genefilter 1.13.7, geneplotter 1.13.5, hgu95av2cdf 1.15.0, lattice 0.14-16, matchprobes 1.7.4, simpleaffy 2.9.1, survival 2.30, xtable 1.4-2