# Association of Genomic Features with Integration

Charles C. Berry

September 10, 2007

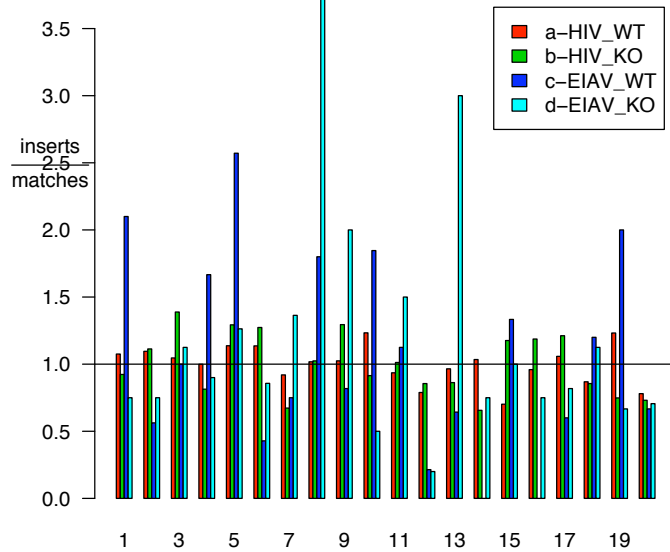## Contents

# 1 Introduction

In this document, I examine the association of integration sites with various genomic features.

The data consist of both actual integration sites and sets of control sites, each set chosen to match the spacing (in bases) from the nearest restriction site (according to the direction in which the sequence was read) to an integration site. The numbers of insertion and matching sites for several data sets are shown below:

```
                 type
Origin.of.data.set insertion match
        a-HIV_WT      1105  3315
        b-HIV_KO       496  1482
        c-EIAV_WT       70   210
        d-EIAV_KO       86   255
```

The advantage of choosing 'control' sites that match the spacing from the nearest restriction site is that biases due to location and density of restriction sites are eliminated by applying the classical multinomial logit model (reviewed in [2]). This model allows regression procedures to be applied to the study of integration intensity as a function of genomic features. The `clogit` function of the R `survival` library) implements estimation and fitting for such models along with the usual likelihood ratio and Wald tests.

The distribution of relative frequency of insertions across the chromosomes is given in this barplot:

It seems evident that there are some chromosomes that are particularly favored for integration. This is reinforced by a test of statistical significance. The test performed used the likelihood ratio statistic for the multinomial logit model (reviewed in [2]) as implemented by the `clogit` function of the R `survival` library). The null hypothesis tested is that the ratio of true integration events to matched control sites is constant across all chromosomes. This test attains a p-value of 0.13482.

# 2 Preference for Genes

## 2.1 refGene Genes

Here we examine the preference that integration events have for genes. In the following plot we show the relative frequency of integrations in genes according to the 'refGene' annotation. The bars grouped over the label "In Gene" give the relative frequency of integration events (compared to control sites) between bases located within refGene gene annotations, while the label "Not in Gene" give the relative frequency of integration events (compared to control sites) between bases not located within refGene gene annotations.
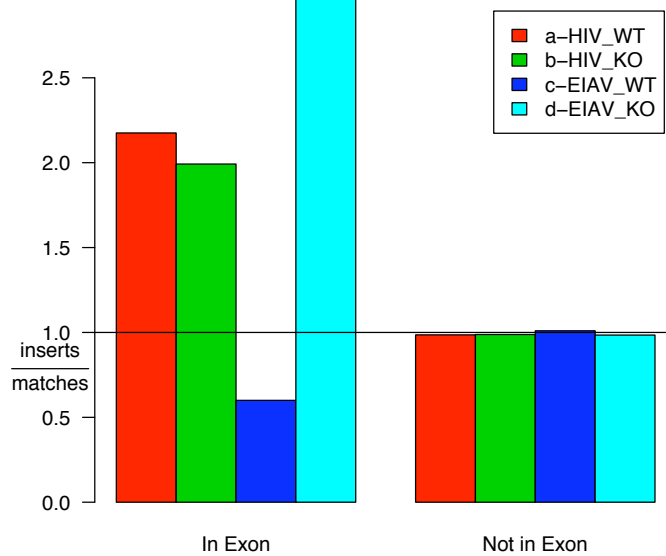


It seems evident that there is a strong tendency for insertions to occur in genes. A formal test of significance bears this out with a p-value of $< 2.22e - 16$. Also, it appears that the tendency of different viruses to integrate into genes varies, and a test for this hypothesis attains $1.2621e - 05$. Here is the table of coefficients of the log ratio of intensities for true insertion sites versus control insertion sites along with their standard errors, z statistics, and p-values for each data set:

```
            coef      se       z         p
a-HIV_WT   1.090  0.0744  14.600  3.51e-48
b-HIV_KO   0.520  0.1110   4.700  2.55e-06
c-EIAV_WT  1.090  0.2810   3.890  1.01e-04
```

```
d-EIAV_KO 0.219 0.2580  0.847 3.97e-01
```

As is evident, there are some differences in the coefficients. The largest coefficient is seen in the c-EIAV$_W T dataset, while the smallest is seen in the d-EIAV_KO dataset.$

In the following plot we show the relative frequency of insertions in exons according to the 'refGene' annotation. The bars grouped over the label "In Exon" give the relative frequency of integration events (compared to control sites) between bases located in exons according to the Acembly annotation, while the label "Not in Exon" give the relative frequency of integration events (compared to control sites) between bases not located in exons according to the Acembly gene annotation.



Here is the table of coefficients of the log ratio of intensities for true insertion sites versus control insertion sites along with their standard errors, z statistics, and p-values for each data set:
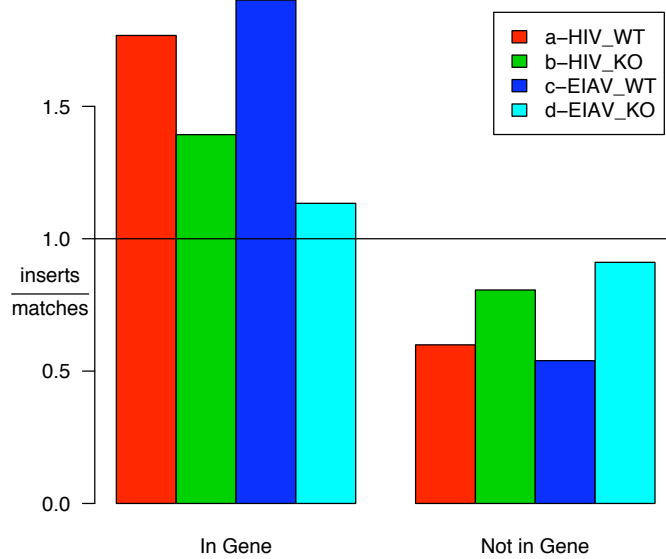
```
            coef    se      z       p
a-HIV_WT   0.129  0.250   0.516  0.606
b-HIV_KO   0.404  0.385   1.050  0.294
c-EIAV_WT -1.060  1.110  -0.954  0.340
d-EIAV_KO  0.970  1.020   0.953  0.341
```

The model on which these coefficients are based include terms for whether the site is in a gene or not. Thus, the effect shown as "In Exon" is net of that

due to being in a gene. Note that in the barplot above the 'Not in Exon' bars include both the introns and intergenic regions, so the impression given by the table may differ from that for the barplot.

## 2.2 ensGenes

Here we examine the preference that insertions have for genes. In the following plot we show the relative frequency of insertions in genes according to the 'ensGene' annotation.
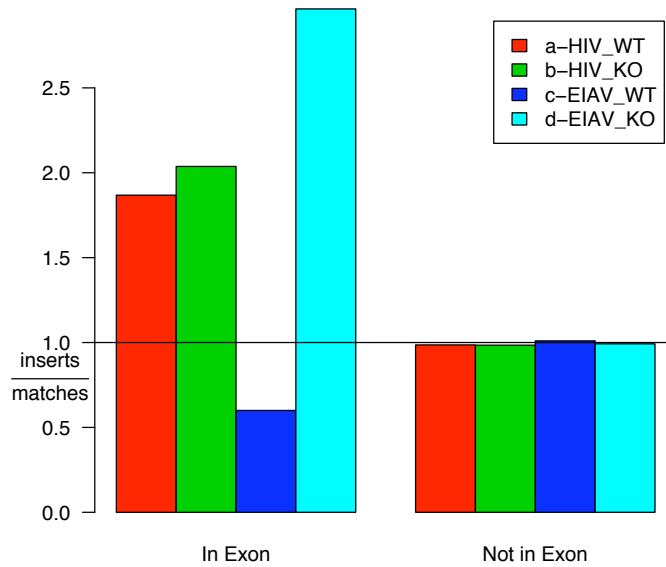


It seems evident that there is a strong tendency for insertions to occur in genes. A formal test of significance bears this out with a p-value of $< 2.22e - 16$. Also, it appears that the tendency of different viruses to integrate into genes varies, and a test for this hypothesis attains $1.3987e - 05$. Here is the table of coefficients of the log ratio of intensities for true insertion sites versus control insertion sites along with their standard errors, z statistics, and p-values for each data set:

```
          coef     se     z         p
a-HIV_WT  1.080 0.0734 14.800 2.55e-49
b-HIV_KO  0.541 0.1060  5.110 3.14e-07
c-EIAV_WT 1.140 0.2760  4.120 3.76e-05
d-EIAV_KO 0.238 0.2520  0.945 3.44e-01
```

7

As is evident, there are some differences in the coefficients. The largest coefficient is seen in the c-EIAV$_W$$Tdataset, while the smallest is seen in the d-EIAV_KO dataset.

In the following plot we show the relative frequency of insertions in exons according to the 'ensGene' annotation.
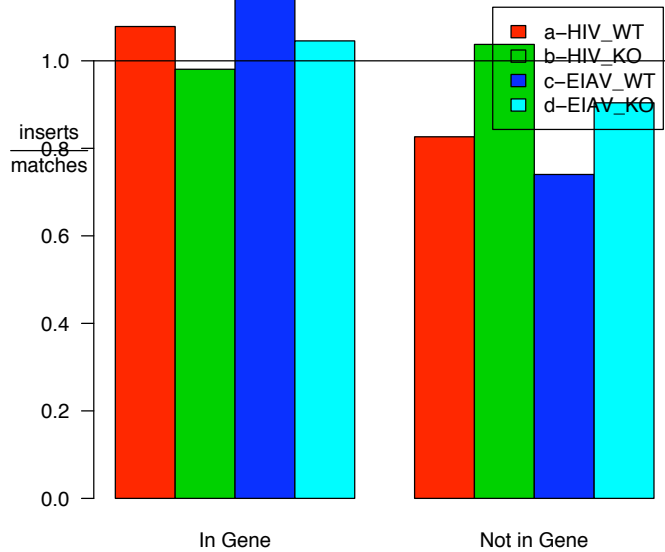


Here is the table of coefficients of the log ratio of intensities for true insertion sites versus control insertion sites along with their standard errors, z statistics, and p-values for each data set:

```
              coef    se      z     p
a-HIV_WT    0.0221 0.228   0.097 0.923
b-HIV_KO    0.4060 0.351   1.160 0.248
c-EIAV_WT  -1.1400 1.110  -1.020 0.307
d-EIAV_KO   0.9920 1.420   0.699 0.485
```

The model on which these coefficients are based include terms for whether the site is in a gene or not. Thus, the effect shown as "In Exon" is net of that due to being in a gene.

## 2.3   genScan Genes

Here we examine the preference that insertions have for genes. In the following plot we show the relative frequency of insertions in genes according to the 'genScan' annotation.
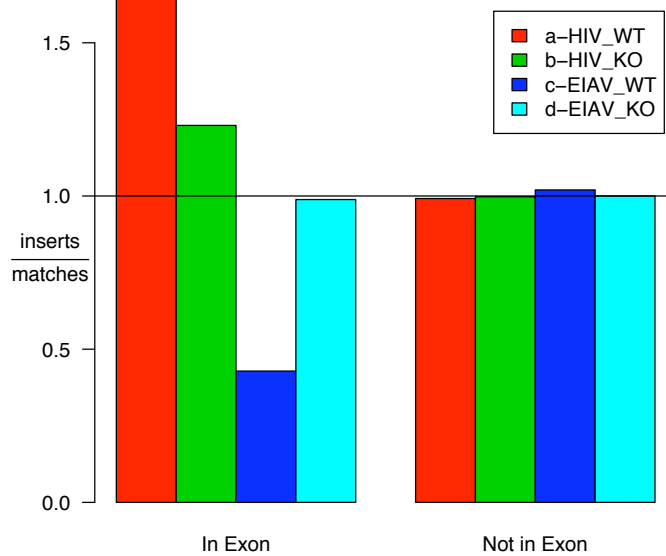
It seems evident that there is a strong tendency for insertions to occur in genes. A formal test of significance bears this out with a p-value of 0.0051861. Also, it appears that the tendency of different viruses to integrate into genes varies, and a test for this hypothesis attains 0.092849. Here is the table of coefficients of the log ratio of intensities for true insertion sites versus control insertion sites along with their standard errors, z statistics, and p-values for each data set:

```
             coef      se       z        p
a-HIV_WT    0.2660  0.0785   3.390  0.000702
b-HIV_KO   -0.0527  0.1080  -0.488  0.626000
c-EIAV_WT   0.4100  0.2940   1.390  0.164000
d-EIAV_KO   0.1210  0.2650   0.456  0.648000
```

As is evident, there are some differences in the coefficients. The largest coefficient is seen in the c-EIAV$_W T dataset, while the smallest is seen in the b-HIV_K O dataset$.

In the following plot we show the relative frequency of insertions in exons according to the 'genScan' annotation.
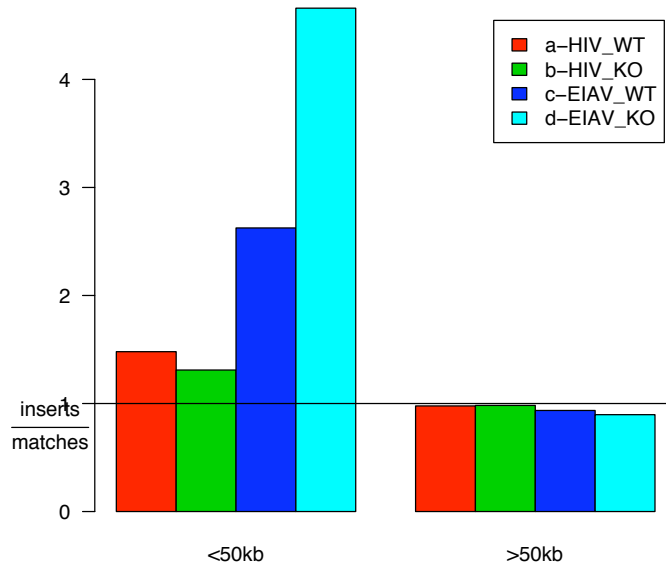
Here is the table of coefficients of the log ratio of intensities for true insertion sites versus control insertion sites along with their standard errors, z statistics, and p-values for each data set:

```
              coef     se        z       p
a-HIV_WT    0.42500 0.262   1.62000 0.105
b-HIV_KO    0.24400 0.459   0.53200 0.595
c-EIAV_WT  -1.02000 1.080  -0.95100 0.342
d-EIAV_KO  -0.00974 1.150  -0.00843 0.993
```

The model on which these coefficients are based include terms for whether the site is in a gene or not. Thus, the effect shown as "In Exon" is net of that due to being in a gene.

## 2.4 oncogenes

Here we examine the preference that insertions have for oncogenes. In the following plot we show the relative frequency of insertions with 50kb of an oncogene 5' end.



A formal test of oncogenic insertion returns p-value of $5.1616e - 05$. The tendency of different viruses to integrate near oncogenes may vary, and a test for this hypothesis attains 0.064614. Here is the table of coefficients of the log ratio of intensities for true insertion sites versus control insertion sites along with their standard errors, z statistics, and p-values for each data set:

```
             coef     se      z        p
a-HIV_WT   -0.418  0.148  -2.82  0.00479
b-HIV_KO   -0.292  0.211  -1.39  0.16600
c-EIAV_WT  -1.100  0.566  -1.94  0.05210
d-EIAV_KO  -1.550  0.483  -3.21  0.00134
a-HIV_WT       NA  0.000     NA       NA
b-HIV_KO       NA  0.000     NA       NA
c-EIAV_WT      NA  0.000     NA       NA
d-EIAV_KO      NA  0.000     NA       NA
```
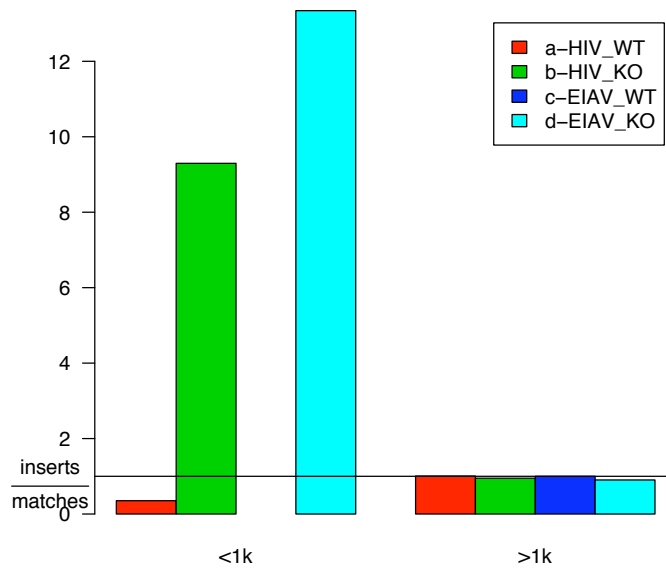
As is evident, there are some differences in the coefficients. The largest coefficient is seen in the b-HIV$_K O dataset, while the smallest is seen in the d-EIAV_K O dataset.$

# 3 CpG Island Neighborhoods

Here we study the effect of being in the neighborhood of CpG Islands. Following Wu et al [3], who found that the neighborhoods within ±1kb of CpG islands are enriched for MLV insertions, we study such neighborhoods.

## 3.1 1 kilobase neighborhoods

The following plot shows the effect of being in or within ±1kb of a CpG island:
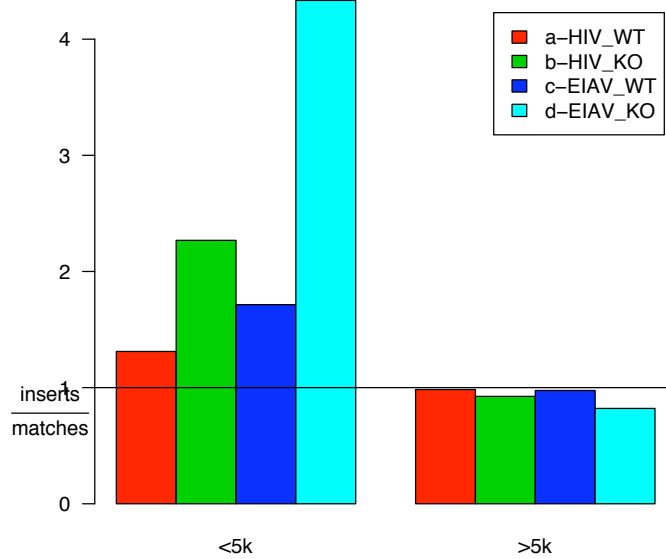


A formal test of significance comparing the difference attains a p-value of $5.2982e - 06$. A test for differences between viruses attains $2.9434e - 10$. Here is the table of coefficients of the log ratio of intensities for true insertion sites versus control insertion sites along with their standard errors, z statistics, and p-values for each data set:

```
          coef    se      z        p
a-HIV_WT  -1.04 0.529  -1.97 4.88e-02
b-HIV_KO   2.23 0.383   5.83 5.57e-09
c-EIAV_WT    NA 0.000     NA       NA
d-EIAV_KO  3.22 1.060   3.05 2.29e-03
```

The largest coefficient is seen in the d-EIAV$_K O dataset, while the smallest is seen in the a-$HIV_W T dataset.

12

## 3.2  5 kilobase neighborhoods

The following plot shows the effect of being in or within ±5kb of a CpG island:
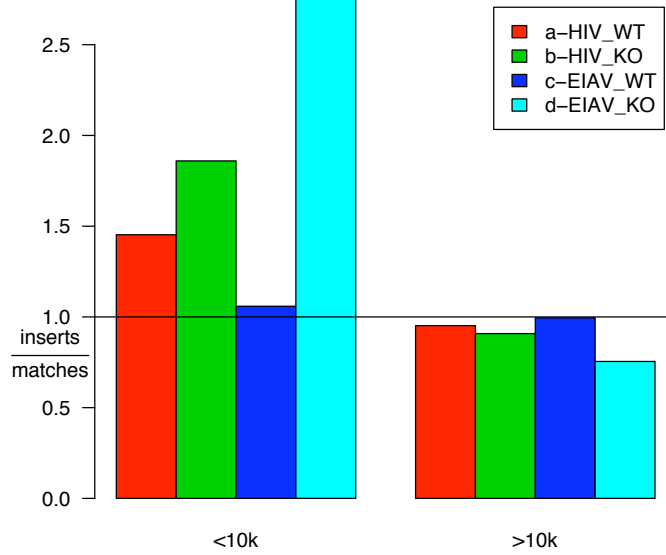


A formal test of significance comparing the difference attains a p-value of $1.5369e - 08$. A test for differences between viruses attains 0.0023532. Here is the table of coefficients of the log ratio of intensities for true insertion sites versus control insertion sites along with their standard errors, z statistics, and p-values for each data set:

```
          coef   se     z        p
a-HIV_WT  0.290 0.146 1.990 4.62e-02
b-HIV_KO  0.893 0.176 5.080 3.81e-07
c-EIAV_WT 0.580 0.654 0.887 3.75e-01
d-EIAV_KO 1.660 0.413 4.020 5.94e-05
```

The largest coefficient is seen in the d-EIAV$_K$$Odataset, while the smallest is seen in the a-$ $HIV_W T dataset$.

## 3.3  10 kilobase neighborhoods

The following plot shows the effect of being in or within ±10kb of a CpG island:
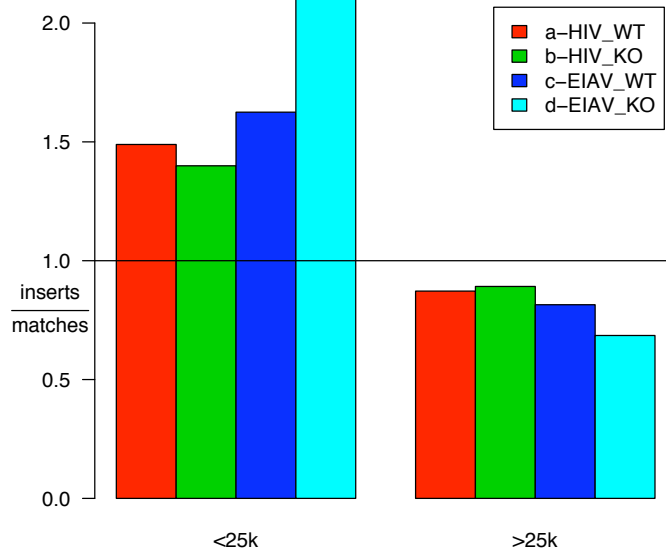
A formal test of significance comparing the difference attains a p-value of $8.6782e-12$. A test for differences between viruses attains $0.018728$. Here is the table of coefficients of the log ratio of intensities for true insertion sites versus control insertion sites along with their standard errors, z statistics, and p-values for each data set:

```
            coef    se     z       p
a-HIV_WT  0.4240 0.105 4.030 5.64e-05
b-HIV_KO  0.7230 0.147 4.900 9.38e-07
c-EIAV_WT 0.0607 0.490 0.124 9.01e-01
d-EIAV_KO 1.3300 0.317 4.180 2.91e-05
```

The largest coefficient is seen in the d-EIAV$_K O dataset, while the smallest is seen in the c-EIAV_W T dataset.$

## 3.4   25 kilobase neighborhoods

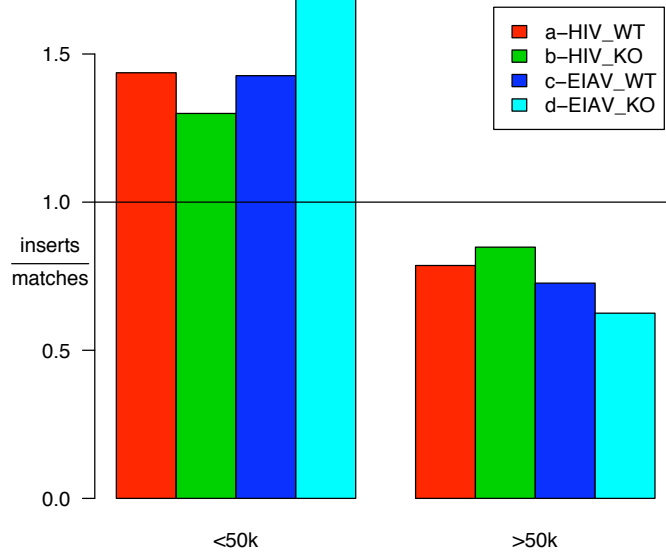The following plot shows the effect of being in or within $\pm 25$kb of a CpG island:

A formal test of significance comparing the difference attains a p-value of $< 2.22e - 16$. A test for differences between viruses attains 0.14831. Here is the table of coefficients of the log ratio of intensities for true insertion sites versus control insertion sites along with their standard errors, z statistics, and p-values for each data set:

```
          coef     se    z         p
a-HIV_WT  0.541  0.0788  6.87  6.58e-12
b-HIV_KO  0.464  0.1180  3.92  8.96e-05
c-EIAV_WT 0.656  0.2910  2.26  2.40e-02
d-EIAV_KO 1.150  0.2800  4.10  4.19e-05
```

The largest coefficient is seen in the d-EIAV$_K$$O dataset, while the smallest is seen in the $b-HIV_K O dataset.$

## 3.5   50 kilobase neighborhoods

The following plot shows the effect of being in or within ±50kb of a CpG island:

A formal test of significance comparing the difference attains a p-value of $< 2.22e - 16$. A test for differences between viruses attains 0.20523. Here is the table of coefficients of the log ratio of intensities for true insertion sites versus control insertion sites along with their standard errors, z statistics, and p-values for each data set:

```
           coef     se     z        p
a-HIV_WT   0.607  0.0714 8.49 2.02e-17
b-HIV_KO   0.430  0.1060 4.06 4.93e-05
c-EIAV_WT  0.671  0.2790 2.40 1.62e-02
d-EIAV_KO  0.974  0.2610 3.72 1.96e-04
```

The largest coefficient is seen in the d-EIAV$_KOdataset, whilethesmallestisseenintheb-$HIV$_KOdataset.$

# 4 Gene Density, Expression 'Density', and CpG Island Density

In this section the association with gene density is examined. For expression analysis, the 'genes' that are counted are the genes represented on the microarray. In addition, we the number of such genes expressed at various levels. The levels are

**low.ex** Count genes whose expression is in the upper half and divide by number of bases

**med.ex** Count genes whose expression is in the upper $1/8^{th}$ and divide by number of bases

**high.ex** Count genes whose expression is in the upper $1/16^{th}$ and divide by number of bases

The bolded terms are used as abbreviations in what follows. The abbreviation **dens** is used to indicate gene density as number of genes per base.

## 4.1 25 kilobase Window

In the barplot that follows we examine the association of insertion sites with gene density in a 25 kilobase window surrounding each locus. More such plots will follow and the method of their construction is always to try to divide the data according to the deciles of density. However, it often happens that there is a very skewed distribution of density and even the $90^{th}$ percentile is zero. In that case, the barplots simply show the sites for which the density is zero and those for which it is non-zero. If there are fewer than ten groups of bars, the groupings contain ten percent of the sites each except for the leftmost grouping which will contain all of the remaining sites.

Also note that the title of the plot contains clues as to its content; the prefix indicates the type of variable studied while the suffix indicates the window width in the number of bases. The p-value given is the result of fitting a cubic polynomial to the gene density values.

The following expression data and probe set were used for this report:

```
[1] "MEF-GSE3400-MGU74Av2"

[1] "MG_U74"

Density data too sparse for barplot

          coef   se    z        p
a-HIV_WT  0.542 0.108 5.03 4.85e-07
b-HIV_KO  0.640 0.159 4.02 5.73e-05
c-EIAV_WT 1.230 0.412 2.98 2.91e-03
d-EIAV_KO 0.753 0.356 2.11 3.46e-02
```

Here are the results for expression density. First, we count just genes that are in the upper half.

```
Density data too sparse for barplot

           coef    se    z        p
a-HIV_WT   0.682 0.134 5.10 3.45e-07
b-HIV_KO   0.818 0.191 4.28 1.84e-05
c-EIAV_WT 1.700 0.508 3.36 7.82e-04
d-EIAV_KO 1.050 0.444 2.36 1.81e-02
```

Now we count genes in the upper $1/8^{th}$:
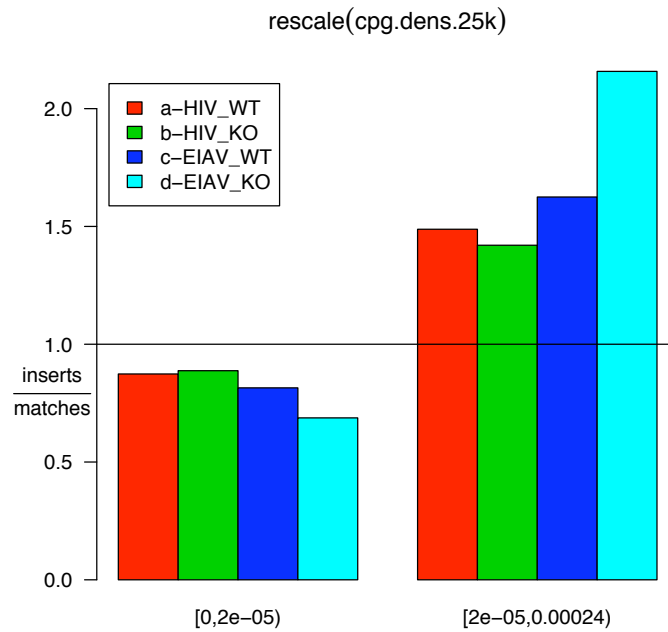
```
Density data too sparse for barplot

          coef    se     z        p
a-HIV_WT  0.845 0.173 4.870 1.10e-06
b-HIV_KO  1.160 0.245 4.720 2.35e-06
c-EIAV_WT 1.500 0.645 2.330 1.98e-02
d-EIAV_KO 0.432 0.635 0.681 4.96e-01
```

And here we count genes in the upper $1/16^{th}$:

```
Density data too sparse for barplot

           coef    se     z        p
a-HIV_WT   0.978 0.246 3.980 6.77e-05
b-HIV_KO   0.973 0.354 2.750 6.00e-03
c-EIAV_WT  1.790 0.866 2.070 3.86e-02
d-EIAV_KO  0.693 0.913 0.759 4.48e-01
```

Here the effect of density of CpG islands is studied:



rescale(cpg.dens.25k)

```
          coef    se    z        p
a-HIV_WT  0.539 0.079 6.82 8.94e-12
b-HIV_KO  0.478 0.119 4.03 5.57e-05
c-EIAV_WT 0.656 0.291 2.26 2.40e-02
d-EIAV_KO 1.190 0.284 4.18 2.89e-05
```

## 4.2   50 kilobase Window

In the barplot that follows we examine the association of insertion sites with expression density in a 50 kilobase window surrounding each locus. First, we count just the number of genes represented on the chip.

rescale(dens.50k)

```
            coef      se     z         p
a-HIV_WT  0.581  0.0864  6.72  1.83e-11
b-HIV_KO  0.738  0.1300  5.66  1.50e-08
c-EIAV_WT 0.761  0.3360  2.26  2.36e-02
d-EIAV_KO 0.923  0.2860  3.23  1.23e-03
```

Here are the results for expression density.  First, we count just genes that are in the upper half.

rescale(low.ex.50k)



```
           coef    se     z        p
a-HIV_WT   0.693 0.102 6.79 1.10e-11
b-HIV_KO   0.889 0.153 5.81 6.41e-09
c-EIAV_WT  1.620 0.445 3.64 2.78e-04
d-EIAV_KO  1.380 0.370 3.73 1.94e-04
```

Now we count genes in the upper $1/8^{th}$:

```
Density data too sparse for barplot

            coef     se     z        p
a-HIV_WT   0.914  0.132  6.92  4.39e-12
b-HIV_KO   1.240  0.194  6.39  1.67e-10
c-EIAV_WT  1.390  0.540  2.57  1.03e-02
d-EIAV_KO  0.943  0.455  2.07  3.82e-02
```

And here we count genes in the upper $1/16^{th}$:

```
Density data too sparse for barplot

            coef    se    z         p
a-HIV_WT   0.896 0.177 5.06 4.15e-07
b-HIV_KO   1.130 0.259 4.37 1.27e-05
c-EIAV_WT  1.610 0.730 2.20 2.75e-02
d-EIAV_KO  0.916 0.606 1.51 1.30e-01
```

Here the effect of density of CpG islands is studied:



rescale(cpg.dens.50k)

```
            coef     se    z        p
a-HIV_WT   0.611  0.0714  8.55  1.21e-17
b-HIV_KO   0.422  0.1060  3.98  6.82e-05
c-EIAV_WT  0.671  0.2790  2.40  1.62e-02
d-EIAV_KO  0.974  0.2610  3.72  1.96e-04
```

## 4.3   100 kilobase Window

In the barplot that follows we examine the association of insertion sites with expression density in a 100 kilobase window surrounding each locus. First, we count just the number of genes represented on the chip.
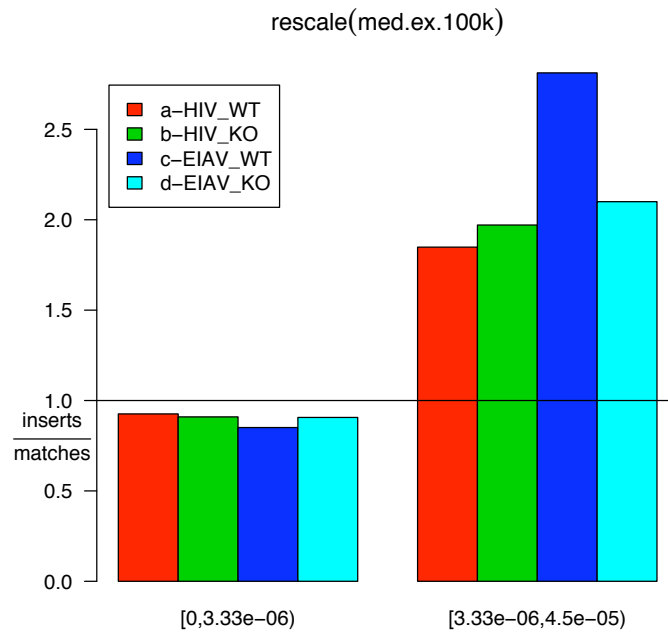
rescale(dens.100k)

|  | coef | se | z | p |
|---|---|---|---|---|
| a-HIV_WT | 0.510 | 0.0749 | 6.81 | 9.51e-12 |
| b-HIV_KO | 0.571 | 0.1140 | 5.03 | 4.99e-07 |
| c-EIAV_WT | 1.170 | 0.3100 | 3.75 | 1.74e-04 |
| d-EIAV_KO | 0.912 | 0.2570 | 3.55 | 3.78e-04 |

Here are the results for expression density. First, we count just genes that are in the upper half.

rescale(low.ex.100k)



```
           coef      se    z        p
a-HIV_WT   0.645  0.0849 7.59 3.08e-14
b-HIV_KO   0.664  0.1250 5.30 1.19e-07
c-EIAV_WT  1.280  0.3480 3.67 2.44e-04
d-EIAV_KO  1.230  0.3080 3.98 6.99e-05
```

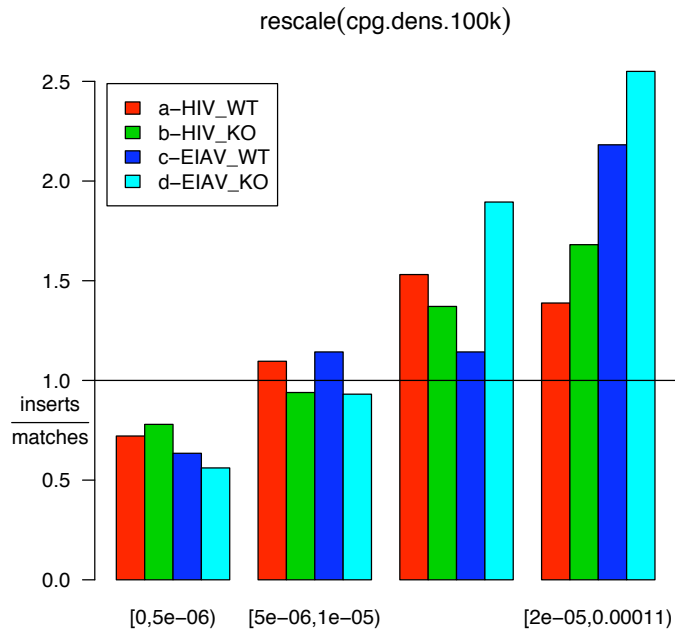Now we count genes in the upper $1/8^{th}$:



rescale(med.ex.100k)

```
          coef   se    z       p
a-HIV_WT  0.721 0.106 6.81 9.50e-12
b-HIV_KO  0.810 0.151 5.36 8.21e-08
c-EIAV_WT 1.270 0.417 3.04 2.33e-03
d-EIAV_KO 0.823 0.372 2.21 2.69e-02
```

And here we count genes in the upper $1/16^{th}$:

```
Density data too sparse for barplot

          coef    se    z         p
a-HIV_WT  0.587 0.140 4.20 2.62e-05
b-HIV_KO  1.000 0.198 5.04 4.61e-07
c-EIAV_WT 0.981 0.486 2.02 4.35e-02
d-EIAV_KO 0.965 0.518 1.86 6.22e-02
```

Here the effect of density of CpG islands is studied:



rescale(cpg.dens.100k)

|          | coef  | se     | z    | p        |
|----------|-------|--------|------|----------|
| a-HIV_WT | 0.570 | 0.0727 | 7.85 | 4.30e-15 |
| b-HIV_KO | 0.579 | 0.1090 | 5.31 | 1.11e-07 |
| c-EIAV_WT| 0.782 | 0.3010 | 2.60 | 9.35e-03 |
| d-EIAV_KO| 1.140 | 0.2690 | 4.24 | 2.21e-05 |

## 4.4   250 kilobase Window

In the barplot that follows we examine the association of insertion sites with expression density in a 250 kilobase window surrounding each locus. First, we count just the number of genes represented on the chip.

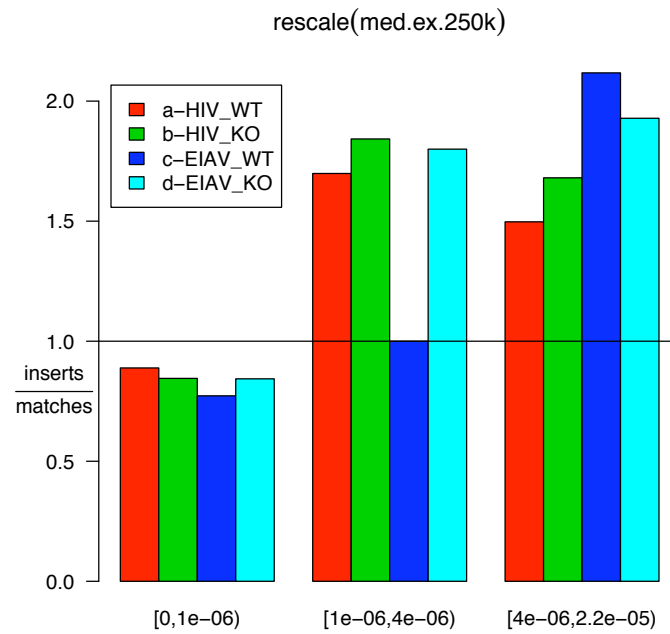rescale(dens.250k)

```
           coef     se    z        p
a-HIV_WT  0.501 0.0705 7.11 1.19e-12
b-HIV_KO  0.471 0.1050 4.49 7.18e-06
c-EIAV_WT 0.705 0.2790 2.53 1.15e-02
d-EIAV_KO 0.602 0.2540 2.37 1.77e-02
```

Here are the results for expression density. First, we count just genes that are in the upper half.
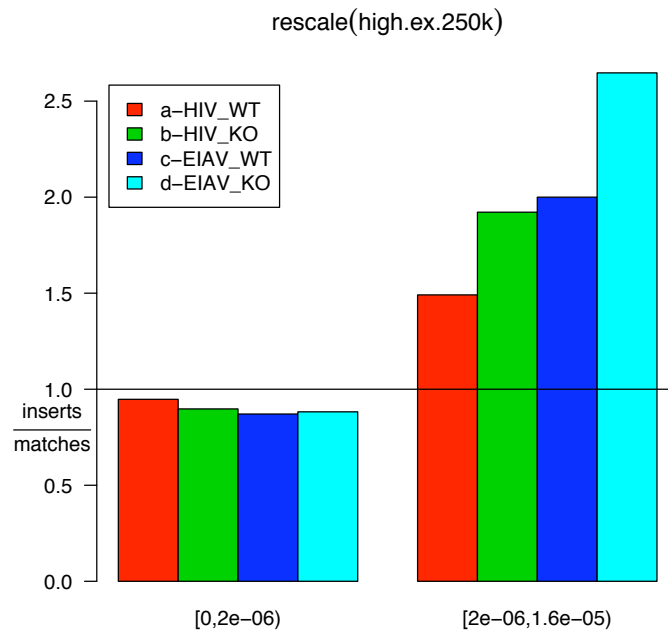


rescale(low.ex.250k)

```
          coef     se    z         p
a-HIV_WT  0.550  0.0719  7.64  2.15e-14
b-HIV_KO  0.546  0.1080  5.06  4.23e-07
c-EIAV_WT 0.856  0.2910  2.94  3.32e-03
d-EIAV_KO 0.897  0.2710  3.31  9.37e-04
```
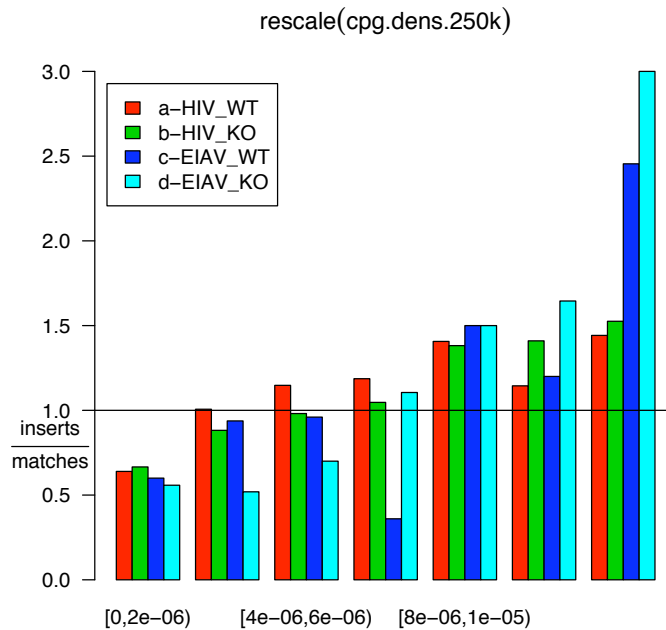
Now we count genes in the upper $1/8^{th}$:

**rescale(med.ex.250k)**



|          | coef  | se     | z    | p        |
|----------|-------|--------|------|----------|
| a-HIV_WT | 0.541 | 0.0813 | 6.65 | 2.85e-11 |
| b-HIV_KO | 0.717 | 0.1210 | 5.91 | 3.43e-09 |
| c-EIAV_WT| 0.981 | 0.3250 | 3.01 | 2.57e-03 |
| d-EIAV_KO| 0.841 | 0.3080 | 2.73 | 6.28e-03 |

And here we count genes in the upper $1/16^{th}$:



rescale(high.ex.250k)

```
          coef    se    z       p
a-HIV_WT  0.466 0.100 4.65 3.31e-06
b-HIV_KO  0.779 0.142 5.50 3.85e-08
c-EIAV_WT 0.764 0.365 2.10 3.62e-02
d-EIAV_KO 1.000 0.380 2.64 8.22e-03
```

Here the effect of density of CpG islands is studied:



rescale(cpg.dens.250k)

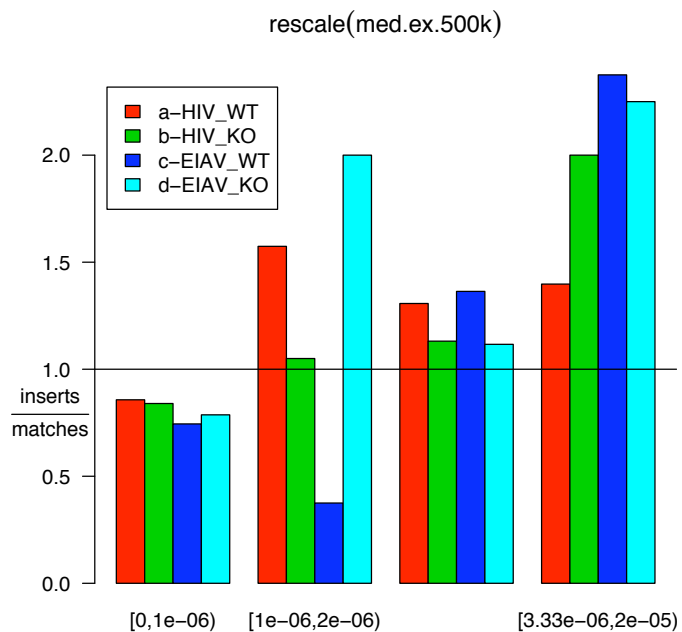```
           coef     se    z        p
a-HIV_WT  0.428  0.0701  6.11  1.01e-09
b-HIV_KO  0.527  0.1050  5.04  4.68e-07
c-EIAV_WT 0.558  0.2810  1.99  4.69e-02
d-EIAV_KO 1.110  0.2580  4.29  1.82e-05
```

## 4.5   500 kilobase Window

In the barplot that follows we examine the association of insertion sites with expression density in a 500 kilobase window surrounding each locus. First, we count just the number of genes represented on the chip.

    @ ■count-500k,echo=F,fig=T■=
    dens.var <- expression(rescale(dens.500k))
    ■gene-dense■
    @

```
           coef     se    z        p
a-HIV_WT  0.428  0.0701  6.11  1.01e-09
b-HIV_KO  0.527  0.1050  5.04  4.68e-07
c-EIAV_WT 0.558  0.2810  1.99  4.69e-02
d-EIAV_KO 1.110  0.2580  4.29  1.82e-05
```

Here are the results for expression density. First, we count just genes that are in the upper half.
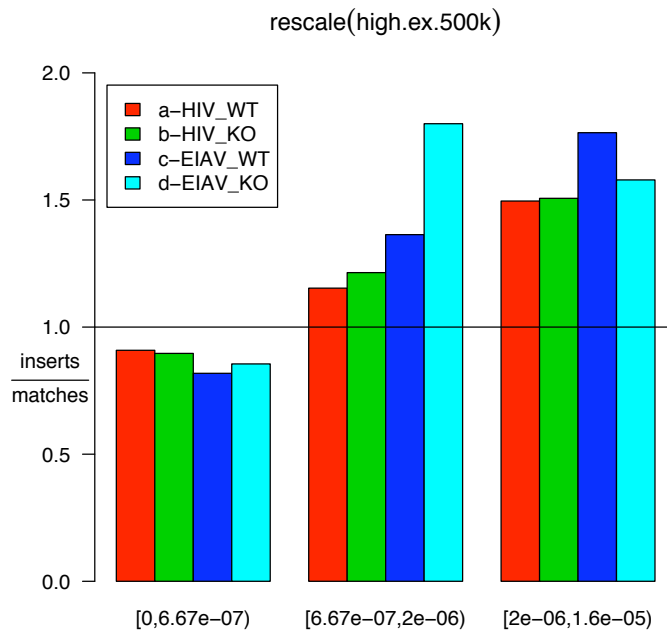


rescale(low.ex.500k)

|          | coef  | se     | z    | p        |
| -------- | ----- | ------ | ---- | -------- |
| a-HIV_WT | 0.487 | 0.0697 | 6.98 | 2.88e-12 |
| b-HIV_KO | 0.534 | 0.1050 | 5.09 | 3.55e-07 |
| c-EIAV_WT | 0.624 | 0.2930 | 2.13 | 3.30e-02 |
| d-EIAV_KO | 0.959 | 0.2650 | 3.62 | 2.96e-04 |

Now we count genes in the upper $1/8^{th}$:



rescale(med.ex.500k)

```
          coef     se    z        p
a-HIV_WT  0.438 0.0719 6.09 1.14e-09
b-HIV_KO  0.539 0.1080 5.00 5.70e-07
c-EIAV_WT 0.689 0.2930 2.36 1.85e-02
d-EIAV_KO 0.645 0.2610 2.47 1.34e-02
```

And here we count genes in the upper $1/16^{th}$:



rescale(high.ex.500k)

```
          coef     se    z        p
a-HIV_WT  0.441 0.0829 5.32 1.05e-07
b-HIV_KO  0.509 0.1210 4.21 2.52e-05
c-EIAV_WT 0.672 0.3000 2.24 2.51e-02
d-EIAV_KO 0.703 0.3010 2.34 1.95e-02
```

Here the effect of density of CpG islands is studied:

rescale(cpg.dens.500k)



```
           coef     se    z        p
a-HIV_WT  0.275  0.0695  3.95  7.85e-05
b-HIV_KO  0.557  0.1050  5.29  1.24e-07
c-EIAV_WT 0.516  0.2820  1.83  6.75e-02
d-EIAV_KO 0.994  0.2620  3.79  1.48e-04
```

## 4.6  1 megabase Window

In the barplot that follows we examine the association of insertion sites with
expression density in a 1 megabase window surrounding each locus. First, we
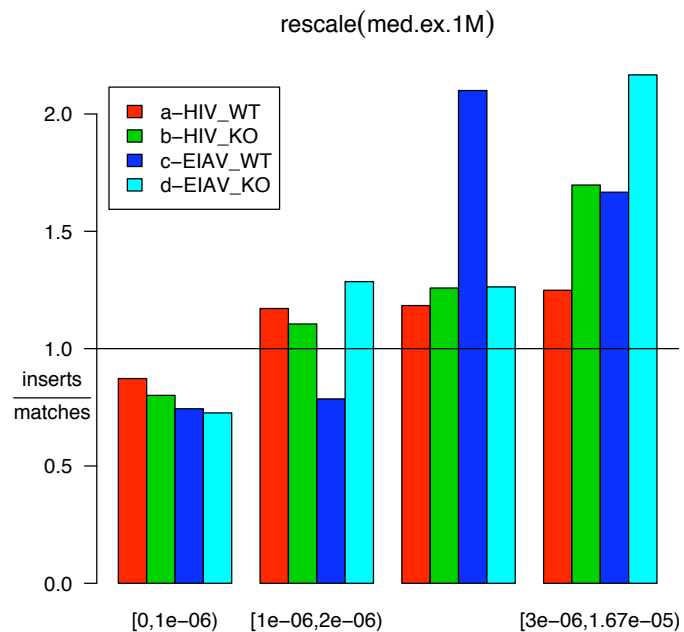count just the number of genes represented on the chip.

rescale(dens.1M)

|          | coef  | se     | z    | p        |
|----------|-------|--------|------|----------|
| a-HIV_WT | 0.318 | 0.0697 | 4.56 | 5.20e-06 |
| b-HIV_KO | 0.476 | 0.1040 | 4.59 | 4.48e-06 |
| c-EIAV_WT | 0.528 | 0.2860 | 1.85 | 6.50e-02 |
| d-EIAV_KO | 0.491 | 0.2550 | 1.92 | 5.47e-02 |

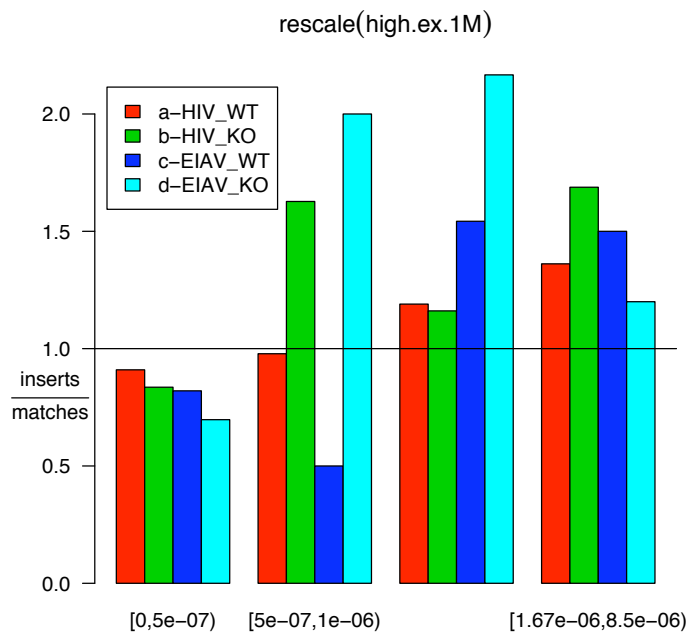Here are the results for expression density. First, we count just genes that are in the upper half.

rescale(low.ex.1M)



```
          coef    se     z       p
a-HIV_WT  0.336  0.0697  4.83  1.39e-06
b-HIV_KO  0.488  0.1030  4.74  2.11e-06
c-EIAV_WT 0.611  0.2850  2.14  3.21e-02
d-EIAV_KO 0.780  0.2550  3.06  2.20e-03
```
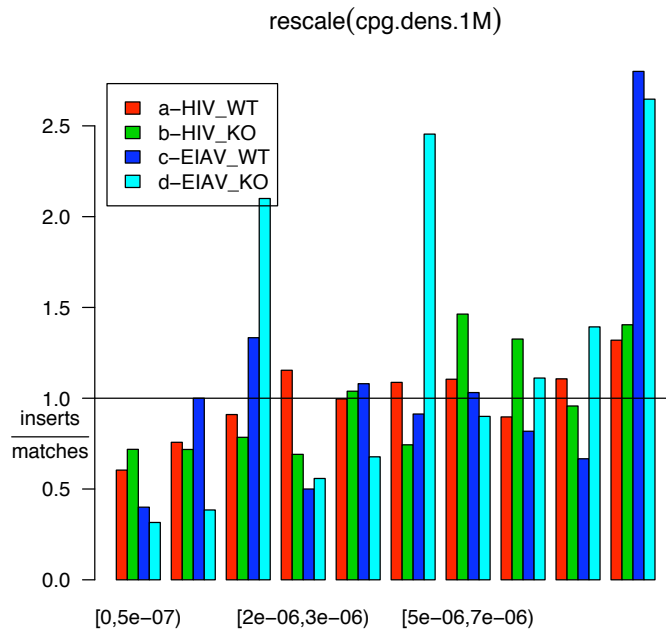
Now we count genes in the upper $1/8^{th}$:

### rescale(med.ex.1M)



```
          coef      se    z         p
a-HIV_WT  0.311  0.0695  4.48  7.57e-06
b-HIV_KO  0.525  0.1050  5.00  5.85e-07
c-EIAV_WT 0.401  0.2860  1.40  1.61e-01
d-EIAV_KO 0.744  0.2600  2.87  4.15e-03
```

And here we count genes in the upper $1/16^{th}$:



rescale(high.ex.1M)

```
          coef      se    z        p
a-HIV_WT  0.251  0.0724  3.48  5.10e-04
b-HIV_KO  0.513  0.1060  4.84  1.33e-06
c-EIAV_WT 0.410  0.2780  1.47  1.40e-01
d-EIAV_KO 0.989  0.2780  3.55  3.82e-04
```
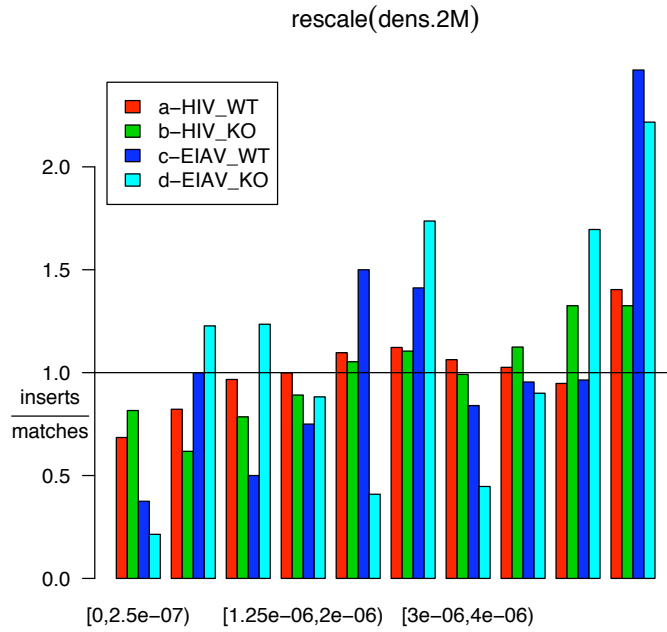
Here the effect of density of CpG islands is studied:



rescale(cpg.dens.1M)

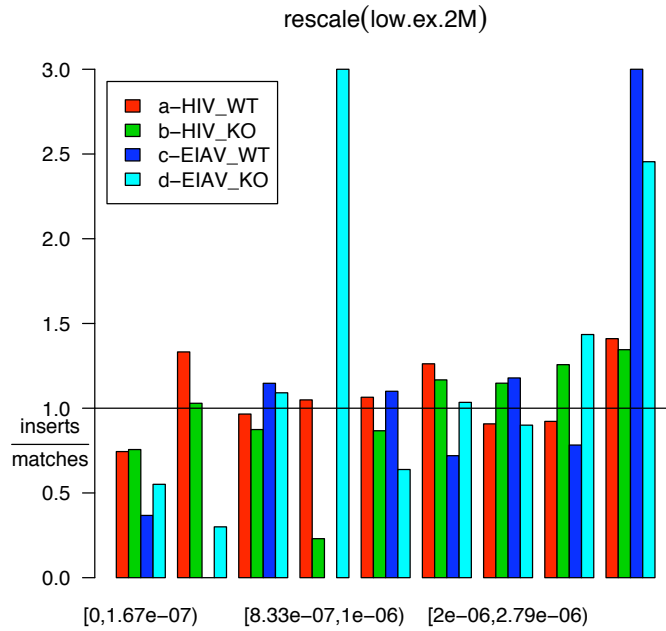|          | coef  | se     | z    | p        |
|----------|-------|--------|------|----------|
| a-HIV_WT | 0.205 | 0.0697 | 2.94 | 0.003250 |
| b-HIV_KO | 0.381 | 0.1050 | 3.64 | 0.000276 |
| c-EIAV_WT | 0.512 | 0.2870 | 1.78 | 0.075000 |
| d-EIAV_KO | 0.692 | 0.2580 | 2.68 | 0.007340 |

## 4.7  2 megabase Window

In the barplot that follows we examine the association of insertion sites with expression density in a 2 megabase window surrounding each locus. First, we count just the number of genes represented on the chip.

rescale(dens.2M)



Legend:
- a–HIV_WT
- b–HIV_KO
- c–EIAV_WT
- d–EIAV_KO

x-axis labels: [0,2.5e−07)    [1.25e−06,2e−06)    [3e−06,4e−06)

y-axis: $\frac{inserts}{matches}$
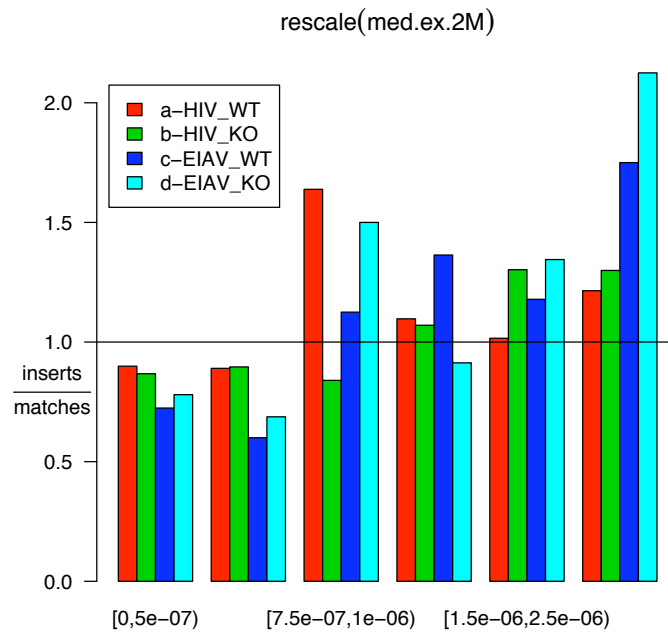
```
          coef      se      z       p
a-HIV_WT  0.175  0.0695  2.510  0.0120
b-HIV_KO  0.336  0.1040  3.240  0.0012
c-EIAV_WT 0.436  0.2780  1.570  0.1170
d-EIAV_KO 0.244  0.2540  0.957  0.3390
```

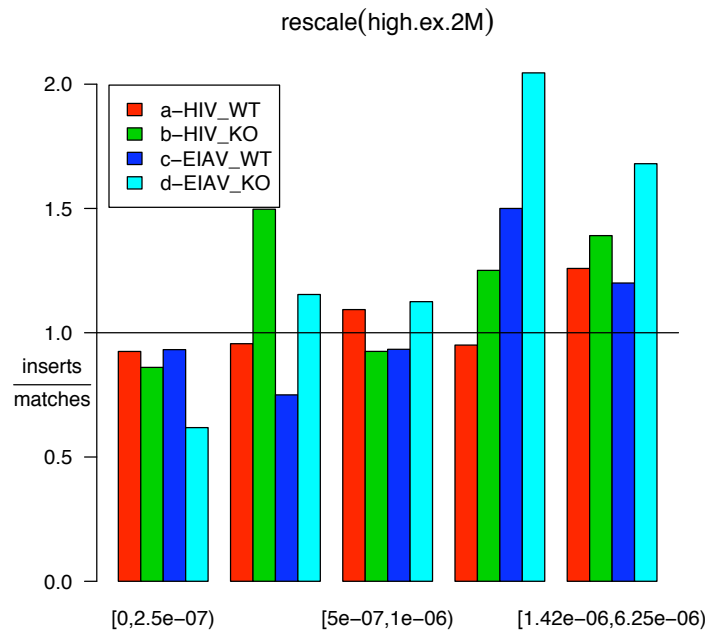Here are the results for expression density. First, we count just genes that are in the upper half.



rescale(low.ex.2M)

```
          coef     se    z          p
a-HIV_WT  0.216  0.0698  3.09  0.002000
b-HIV_KO  0.415  0.1050  3.94  0.000082
c-EIAV_WT 0.531  0.2820  1.88  0.059900
d-EIAV_KO 0.563  0.2540  2.21  0.026800
```

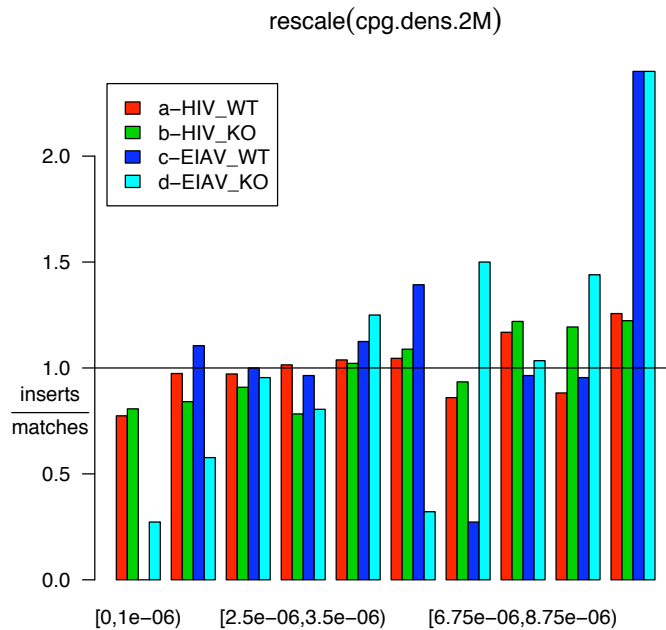Now we count genes in the upper $1/8^{th}$:



rescale(med.ex.2M)

```
           coef      se     z          p
a-HIV_WT  0.241  0.0704  3.43  0.000601
b-HIV_KO  0.324  0.1040  3.13  0.001750
c-EIAV_WT 0.581  0.2850  2.04  0.041200
d-EIAV_KO 0.512  0.2590  1.98  0.048100
```

And here we count genes in the upper $1/16^{th}$:



rescale(high.ex.2M)

```
            coef      se      z        p
a-HIV_WT   0.176  0.0699  2.510  0.01190
b-HIV_KO   0.197  0.1040  1.890  0.05910
c-EIAV_WT  0.144  0.2680  0.537  0.59200
d-EIAV_KO  0.730  0.2650  2.750  0.00589
```
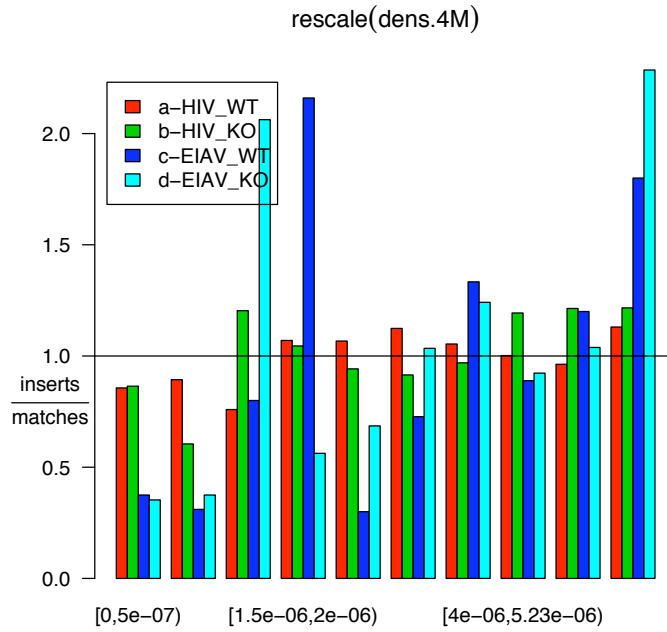
Here the effect of density of CpG islands is studied:



rescale(cpg.dens.2M)

```
          coef     se     z      p
a-HIV_WT  0.116  0.0697  1.670  0.0959
b-HIV_KO  0.257  0.1040  2.480  0.0130
c-EIAV_WT 0.217  0.2820  0.767  0.4430
d-EIAV_KO 0.615  0.2550  2.410  0.0160
```
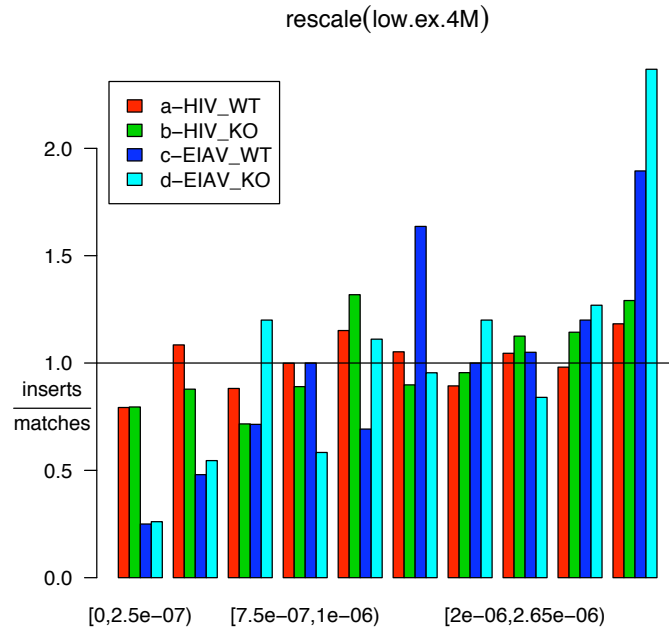
## 4.8   4 megabase Window

In the barplot that follows we examine the association of insertion sites with expression density in a 4 megabase window surrounding each locus. First, we count just the number of genes represented on the chip.
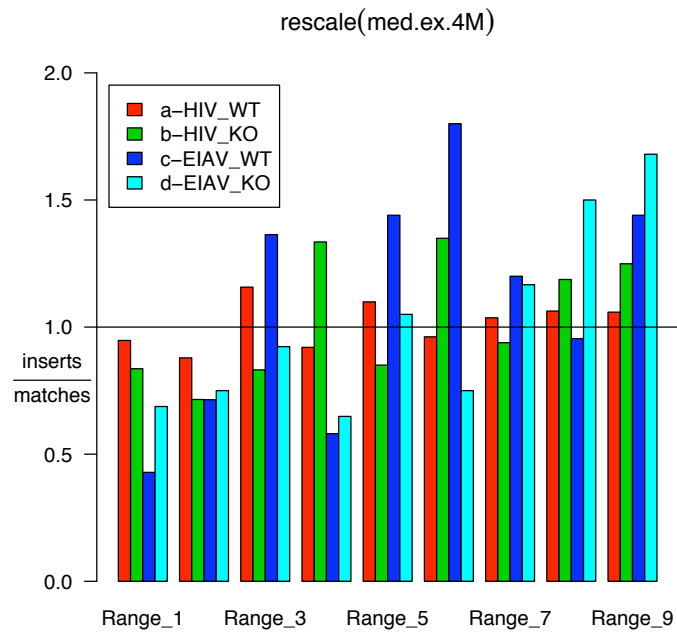
rescale(dens.4M)



|         | coef   | se     | z    | p     |
|---------|--------|--------|------|-------|
| a-HIV_WT | 0.0716 | 0.0695 | 1.03 | 0.303 |
| b-HIV_KO | 0.2020 | 0.1030 | 1.96 | 0.050 |
| c-EIAV_WT | 0.3140 | 0.2740 | 1.15 | 0.252 |
| d-EIAV_KO | 0.5040 | 0.2500 | 2.01 | 0.044 |

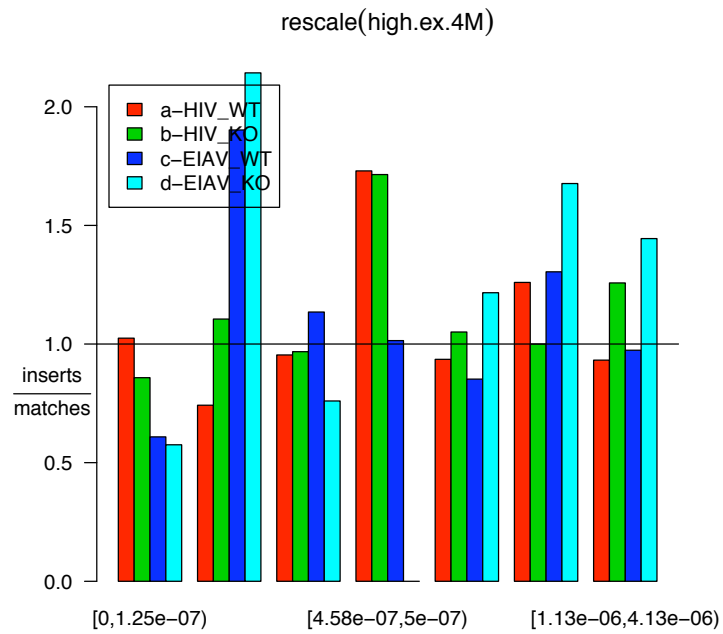Here are the results for expression density. First, we count just genes that are in the upper half.



rescale(low.ex.4M)

|          | coef   | se     | z     | p      |
|----------|--------|--------|-------|--------|
| a-HIV_WT | 0.0395 | 0.0699 | 0.565 | 0.5720 |
| b-HIV_KO | 0.2000 | 0.1020 | 1.960 | 0.0501 |
| c-EIAV_WT| 0.5650 | 0.2690 | 2.100 | 0.0357 |
| d-EIAV_KO| 0.5490 | 0.2530 | 2.170 | 0.0303 |

Now we count genes in the upper $1/8^{th}$:

## rescale(med.ex.4M)



```
           coef     se     z      p
a-HIV_WT  0.0911 0.0699 1.30 0.1920
b-HIV_KO  0.1500 0.1040 1.45 0.1480
c-EIAV_WT 0.4910 0.2740 1.79 0.0739
d-EIAV_KO 0.4010 0.2560 1.57 0.1170
```

And here we count genes in the upper $1/16^{th}$:



rescale(high.ex.4M)

|           | coef   | se     | z     | p      |
|-----------|--------|--------|-------|--------|
| a-HIV_WT  | 0.0669 | 0.0693 | 0.965 | 0.3350 |
| b-HIV_KO  | 0.1680 | 0.1040 | 1.610 | 0.1080 |
| c-EIAV_WT | 0.2130 | 0.2670 | 0.797 | 0.4250 |
| d-EIAV_KO | 0.6110 | 0.2590 | 2.360 | 0.0183 |

Here the effect of density of CpG islands is studied:



rescale(cpg.dens.4M)

|          | coef  | se     | z    | p       |
|----------|-------|--------|------|---------|
| a-HIV_WT | 0.023 | 0.0695 | 0.33 | 0.74100 |
| b-HIV_KO | 0.129 | 0.1040 | 1.24 | 0.21400 |
| c-EIAV_WT| 0.425 | 0.2850 | 1.49 | 0.13600 |
| d-EIAV_KO| 0.724 | 0.2580 | 2.80 | 0.00505 |

## 4.9   8 megabase Window

In the barplot that follows we examine the association of insertion sites with expression density in a 8 megabase window surrounding each locus. First, we count just the number of genes represented on the chip.
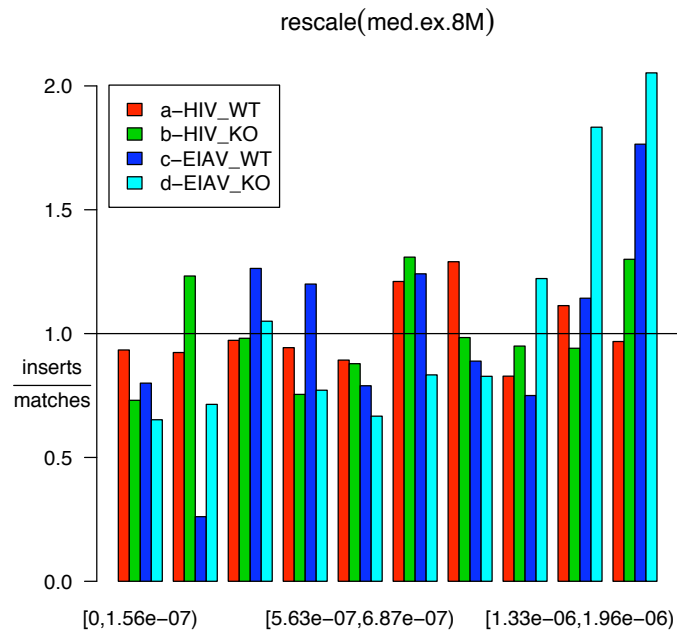
rescale(dens.8M)

|          | coef  | se     | z     | p     |
|----------|-------|--------|-------|-------|
| a-HIV_WT | 0.104 | 0.0696 | 1.490 | 0.137 |
| b-HIV_KO | 0.162 | 0.1050 | 1.550 | 0.122 |
| c-EIAV_WT| 0.203 | 0.2720 | 0.746 | 0.456 |
| d-EIAV_KO| 0.340 | 0.2510 | 1.360 | 0.175 |

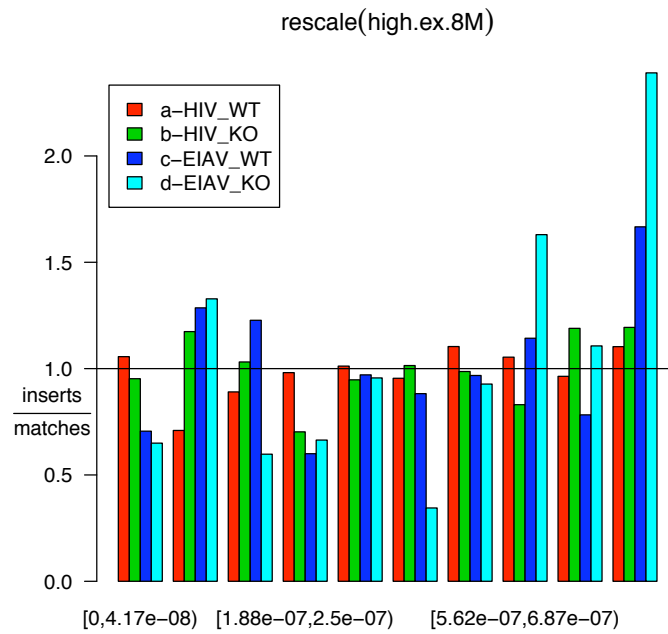Here are the results for expression density. First, we count just genes that are in the upper half.

rescale(low.ex.8M)



```
            coef     se    z       p
a-HIV_WT  0.0844  0.0696  1.21  0.2250
b-HIV_KO  0.1650  0.1050  1.58  0.1140
c-EIAV_WT 0.4330  0.2760  1.57  0.1170
d-EIAV_KO 0.6020  0.2600  2.31  0.0207
```

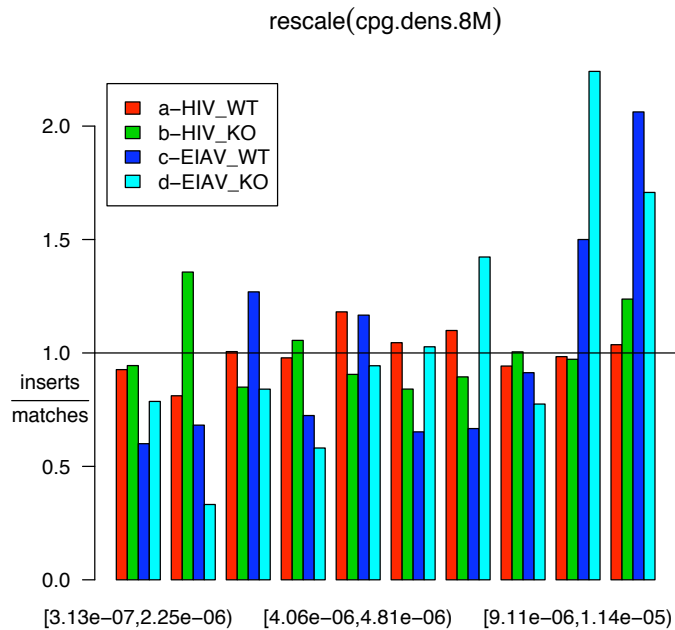57

Now we count genes in the upper $1/8^{th}$:



rescale(med.ex.8M)

|           | coef   | se     | z     | p      |
|-----------|--------|--------|-------|--------|
| a-HIV_WT  | 0.1270 | 0.0698 | 1.830 | 0.0679 |
| b-HIV_KO  | 0.1400 | 0.1030 | 1.360 | 0.1740 |
| c-EIAV_WT | 0.0994 | 0.2820 | 0.353 | 0.7240 |
| d-EIAV_KO | 0.4570 | 0.2590 | 1.760 | 0.0782 |

And here we count genes in the upper $1/16^{th}$:



rescale(high.ex.8M)

```
           coef     se      z       p
a-HIV_WT   0.0875  0.0698  1.250  0.2100
b-HIV_KO   0.1020  0.1040  0.982  0.3260
c-EIAV_WT  0.1270  0.2690  0.471  0.6380
d-EIAV_KO  0.6470  0.2500  2.590  0.0095
```

Here the effect of density of CpG islands is studied:

rescale(cpg.dens.8M)



```
            coef      se      z       p
a-HIV_WT   0.0455  0.0692   0.657  0.51100
b-HIV_KO  -0.0257  0.1050  -0.245  0.80700
c-EIAV_WT  0.2280  0.2750   0.827  0.40800
d-EIAV_KO  0.6780  0.2590   2.610  0.00898
```

## 4.10    16 megabase Window

In the barplot that follows we examine the association of insertion sites with expression density in a 16 megabase window surrounding each locus. First, we count just the number of genes represented on the chip.

rescale(dens.16M)
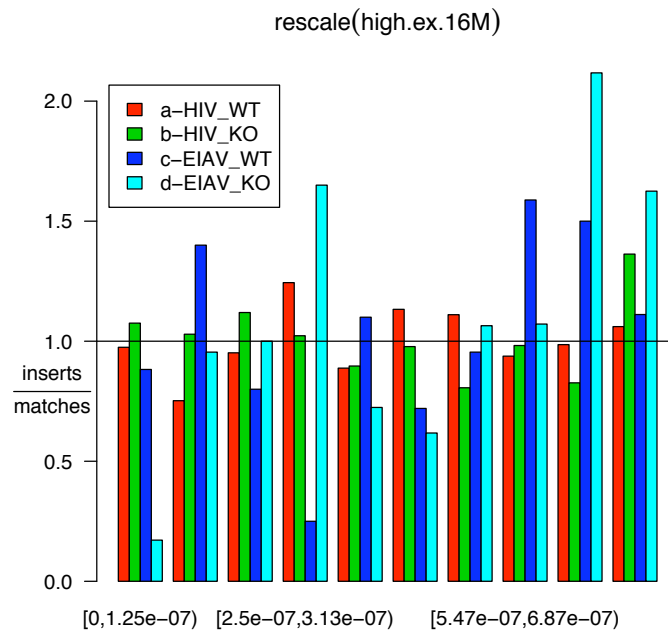


```
            coef      se       z       p
a-HIV_WT  0.0016  0.0693  0.0231  0.982
b-HIV_KO  0.0550  0.1040  0.5280  0.597
c-EIAV_WT 0.4500  0.2760  1.6300  0.103
d-EIAV_KO 0.1850  0.2480  0.7440  0.457
```

Here are the results for expression density. First, we count just genes that are in the upper half.
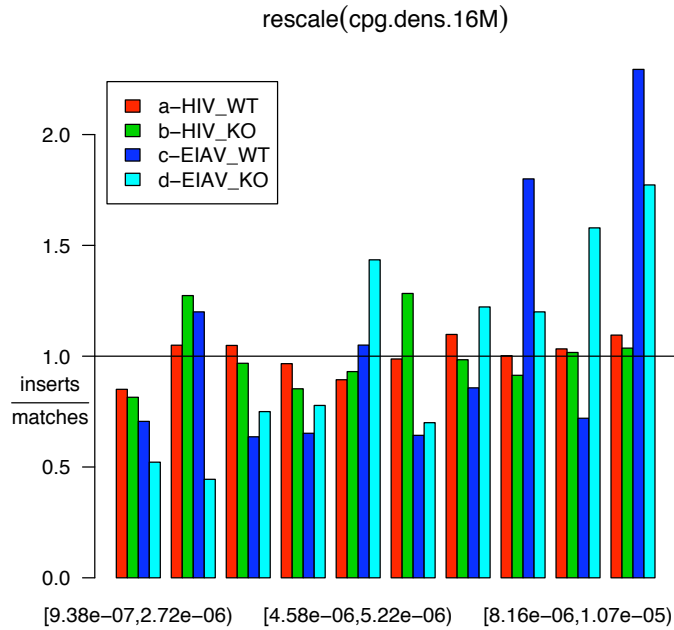
rescale(low.ex.16M)



|          | coef    | se     | z      | p     |
|----------|---------|--------|--------|-------|
| a-HIV_WT | -0.0371 | 0.0692 | -0.537 | 0.592 |
| b-HIV_KO | 0.0339  | 0.1050 | 0.323  | 0.746 |
| c-EIAV_WT | 0.6050 | 0.2840 | 2.130  | 0.033 |
| d-EIAV_KO | 0.3490 | 0.2600 | 1.340  | 0.179 |

Now we count genes in the upper $1/8^{th}$:

## rescale(med.ex.16M)



```
            coef      se        z       p
a-HIV_WT   0.0560  0.0698   0.803  0.422
b-HIV_KO  -0.0155  0.1050  -0.148  0.882
c-EIAV_WT  0.2870  0.2790   1.030  0.303
d-EIAV_KO  0.3810  0.2660   1.430  0.152
```

And here we count genes in the upper $1/16^{th}$:



rescale(high.ex.16M)

```
             coef      se        z      p
a-HIV_WT   0.0977  0.0696   1.400  0.161
b-HIV_KO  -0.0354  0.1040  -0.339  0.734
c-EIAV_WT  0.3360  0.2760   1.220  0.224
d-EIAV_KO  0.3400  0.2630   1.290  0.196
```

Here the effect of density of CpG islands is studied:



rescale(cpg.dens.16M)

```
            coef      se      z       p
a-HIV_WT   0.0898  0.0692  1.300  0.1950
b-HIV_KO   0.0915  0.1040  0.881  0.3780
c-EIAV_WT  0.3920  0.2830  1.390  0.1650
d-EIAV_KO  0.4590  0.2640  1.740  0.0822
```

## 4.11  32 megabase Window

In the barplot that follows we examine the association of insertion sites with expression density in a 32 megabase window surrounding each locus. First, we count just the number of genes represented on the chip.

rescale(dens.32M)

|          | coef   | se     | z     | p     |
|----------|--------|--------|-------|-------|
| a-HIV_WT | 0.0884 | 0.0696 | 1.270 | 0.204 |
| b-HIV_KO | 0.0263 | 0.1040 | 0.252 | 0.801 |
| c-EIAV_WT| 0.4350 | 0.2780 | 1.570 | 0.117 |
| d-EIAV_KO| 0.3470 | 0.2590 | 1.340 | 0.181 |

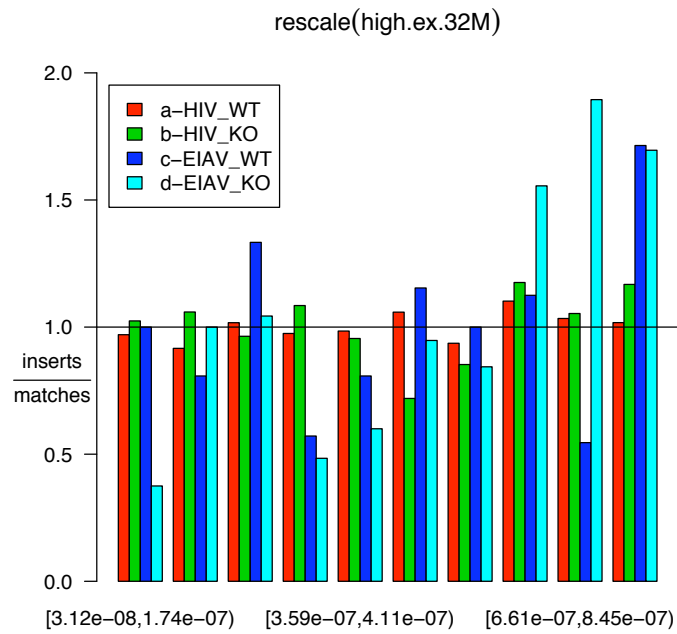Here are the results for expression density. First, we count just genes that are in the upper half.

rescale(low.ex.32M)



```
             coef      se        z       p
a-HIV_WT  -0.00241  0.0694  -0.0347  0.9720
b-HIV_KO   0.04680  0.1040   0.4500  0.6530
c-EIAV_WT  0.55100  0.2910   1.8900  0.0587
d-EIAV_KO  0.59600  0.2540   2.3400  0.0191
```

Now we count genes in the upper $1/8^{th}$:



rescale(med.ex.32M)

|          | coef     | se    | z      | p      |
|----------|----------|-------|--------|--------|
| a-HIV_WT |  0.00833 | 0.069 |  0.121 | 0.9040 |
| b-HIV_KO | -0.01260 | 0.104 | -0.121 | 0.9040 |
| c-EIAV_WT |  0.37000 | 0.288 |  1.280 | 0.2000 |
| d-EIAV_KO |  0.62700 | 0.259 |  2.420 | 0.0156 |

And here we count genes in the upper $1/16^{th}$:



rescale(high.ex.32M)

```
            coef     se       z        p
a-HIV_WT    0.0511 0.069   0.741 0.45900
b-HIV_KO   -0.0229 0.103  -0.223 0.82400
c-EIAV_WT   0.2100 0.290   0.725 0.46900
d-EIAV_KO   0.6800 0.260   2.610 0.00894
```

Here the effect of density of CpG islands is studied:

## rescale(cpg.dens.32M)



```
              coef      se      z      p
a-HIV_WT    0.1020  0.0699  1.460  0.145
b-HIV_KO    0.0699  0.1040  0.674  0.501
c-EIAV_WT   0.2240  0.2860  0.783  0.434
d-EIAV_KO   0.1430  0.2390  0.597  0.550
```

# 5   Juxtaposition with Gene Start and End Positions

## 5.1   Refseq Annotations

In this section we study the effect of juxtaposition in terms of gene start and end positions. The first barplot shows the effect of gene width for those insertions that are located within an RefSeq gene. The table following the barplot shows the p-values for a test of the hypothesis that the proportions in each of the categories that define the bars are equal in the insertions and their matches. This p-value is obtained from the $5 \times 2 \times k$ table of counts defined by gene width category, insertion/match status, and stratum (consisting of an insertion and its matched sites) using a likelihood ratio test for the hypothesis of no association between gene width category and insertion/match status. The test used compared the log-linear model [1] with all two-way configurations to that with no gene width category and insertion/match status configuration.



**refSeq gene.width**

```
 a-HIV_WT  b-HIV_KO c-EIAV_WT d-EIAV_KO
 3.36e-07  2.81e-02  3.00e-01  1.20e-01
```

The next plot uses the width of a non-gene region for insertions that fall into such regions.
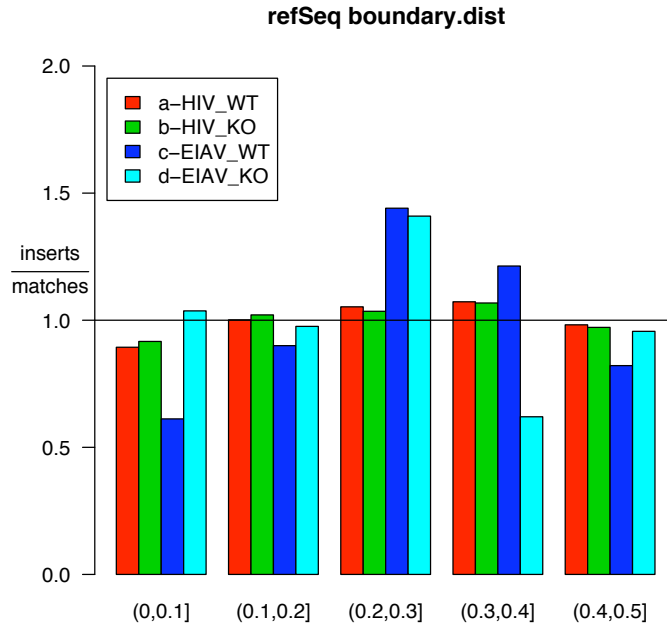
71

**refSeq other.width**

```
 a-HIV_WT  b-HIV_KO c-EIAV_WT d-EIAV_KO
 0.128000  0.000225  0.463000  0.000339
```
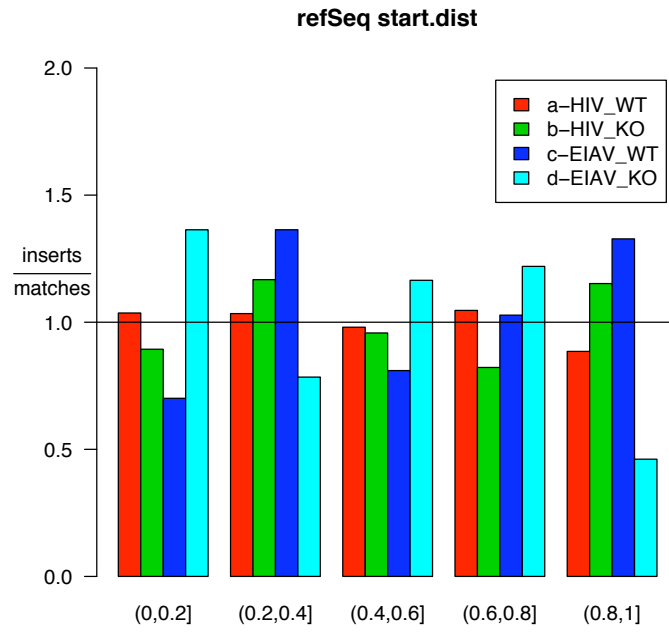
The next plot studies the distance to the nearest boundary between a gene and a non-gene region. The distance is expressed as a fraction of the length of the region. Thus, '0.25' refers to one quarter of the distance from the site to nearest boundary divided by the total width of the region.
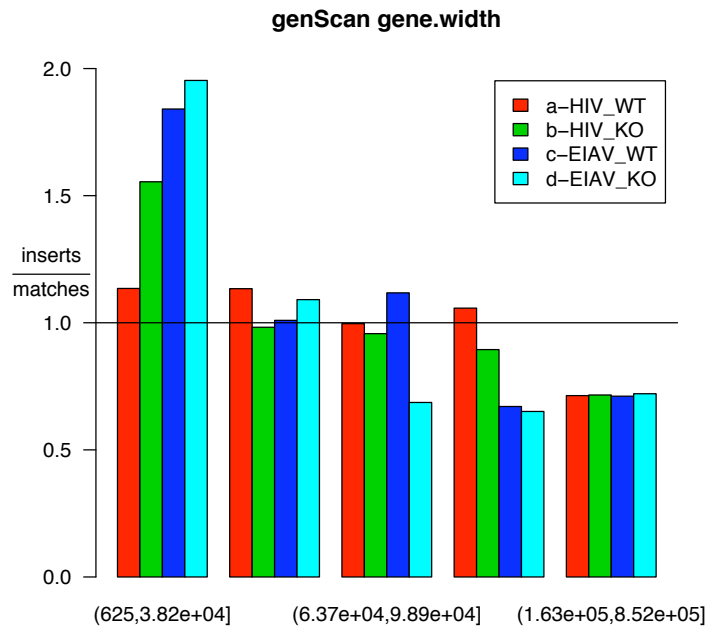
**refSeq boundary.dist**



| a-HIV_WT | b-HIV_KO | c-EIAV_WT | d-EIAV_KO |
|----------|----------|-----------|-----------|
| 0.330    | 0.878    | 0.260     | 0.239     |

This plot studies the effect of nearness to the beginning of a transcript. For sites in genes, it is the distance to the start of the gene divided by the width of the gene. For other sites it is the distance from the site to the nearer gene if that gene boundary is also a transcription starting point. Locations near '0' are relatively near the beginning of transcription, while those near '1' are near the termination of the transcript.

**refSeq start.dist**



| a-HIV_WT | b-HIV_KO | c-EIAV_WT | d-EIAV_KO |
|----------|----------|-----------|-----------|
| 0.468 | 0.206 | 0.117 | 0.168 |

## 5.2   genScan Annotations

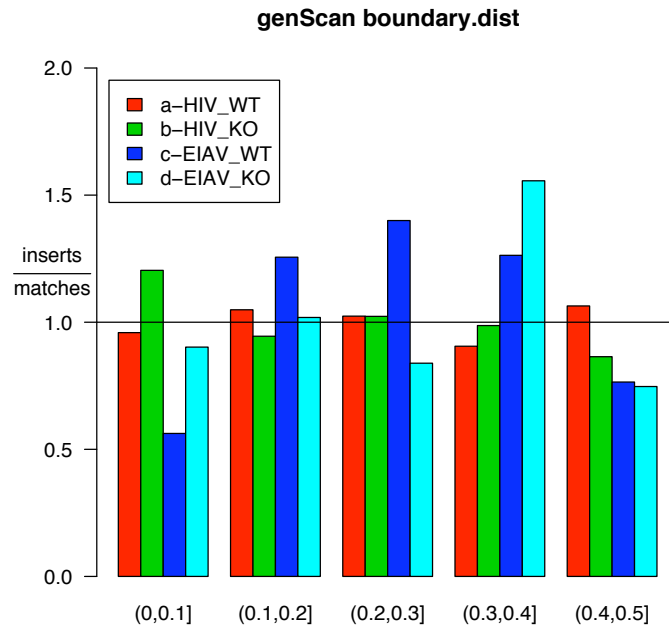**genScan gene.width**



```
a-HIV_WT   b-HIV_KO  c-EIAV_WT  d-EIAV_KO
6.15e-05   3.71e-04   2.92e-01   4.06e-02
```

**genScan other.width**



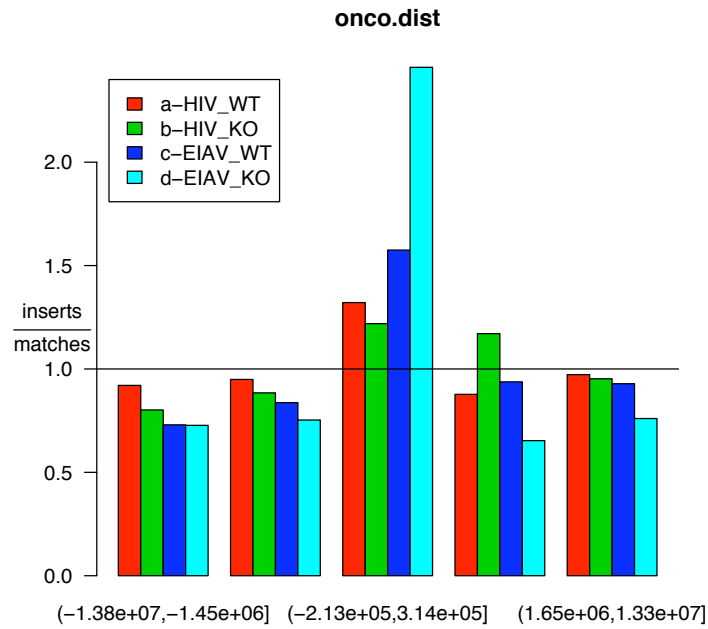| a-HIV_WT | b-HIV_KO | c-EIAV_WT | d-EIAV_KO |
|----------|----------|-----------|-----------|
| 0.0697 | 0.0663 | 0.0564 | 0.4890 |

**genScan boundary.dist**



| a-HIV_WT | b-HIV_KO | c-EIAV_WT | d-EIAV_KO |
|----------|----------|-----------|-----------|
| 0.400    | 0.244    | 0.117     | 0.182     |

**genScan start.dist**



| a-HIV_WT | b-HIV_KO | c-EIAV_WT | d-EIAV_KO |
|----------|----------|-----------|-----------|
| 0.00219  | 0.08270  | 0.08750   | 0.18500   |

# 6   Oncogenes

This plot studies the effect of nearness to the 5' end of an oncogene transcript. Positive values represent distances in which the integration site is upstream of the nearest oncogene 5' end, negative downstream.

**onco.dist**



```
a-HIV_WT  b-HIV_KO c-EIAV_WT d-EIAV_KO
9.40e-05  1.49e-02  2.57e-01  2.96e-05
```

Here is the same plot using absolute distance

**abs abs.onco.dist**

```
a-HIV_WT  b-HIV_KO c-EIAV_WT d-EIAV_KO
3.76e-05  9.31e-02  1.53e-01  3.14e-05
```
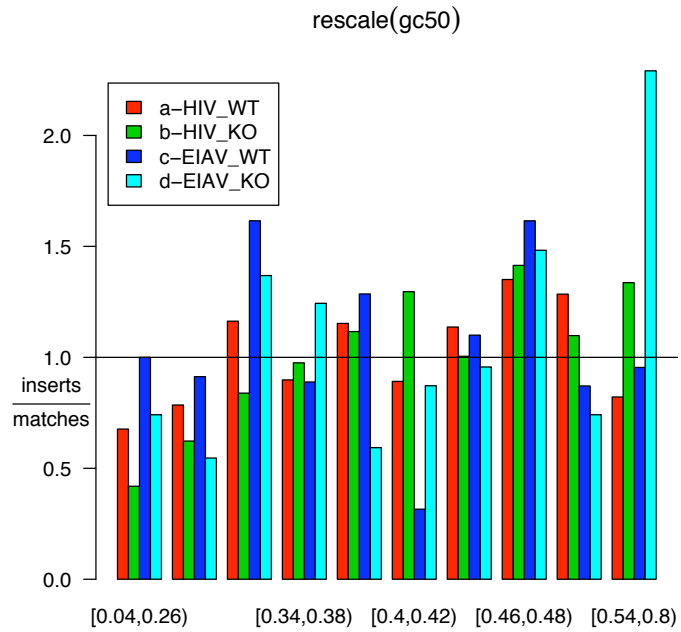
# 7  GC content

Here we study the effect of GC content on insertion. The GC content is taken from the Mouse Genome Draft at GoldenPath from the table

Following the plot is a table of fitted coefficients based on splitting the GC percent data at the median.
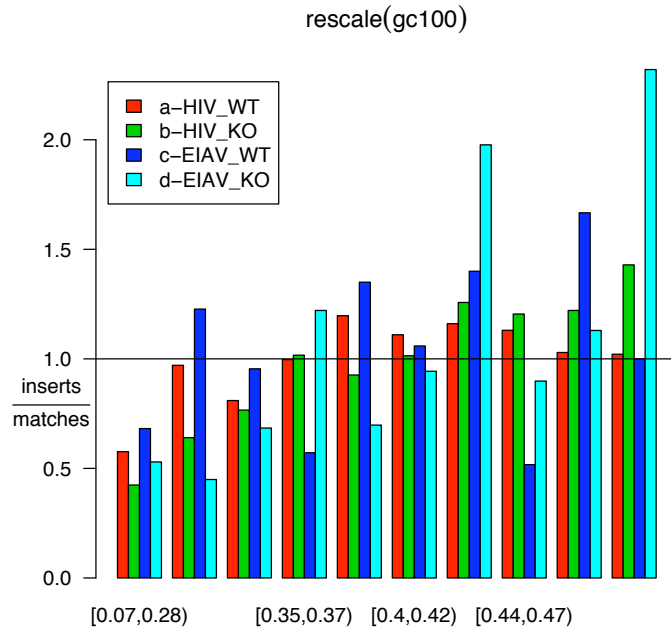


rescale(gc20)

```
              coef      se       z        p
a-HIV_WT    0.242  0.0704   3.440  5.72e-04
b-HIV_KO    0.505  0.1050   4.810  1.51e-06
c-EIAV_WT  -0.411  0.2920  -1.410  1.59e-01
d-EIAV_KO   0.138  0.2610   0.527  5.98e-01

              coef      se       z        p
a-HIV_WT    0.252  0.0718   3.510  0.000448
b-HIV_KO    0.310  0.1040   2.970  0.002930
c-EIAV_WT   0.113  0.2740   0.412  0.680000
d-EIAV_KO   0.418  0.2540   1.640  0.100000
```
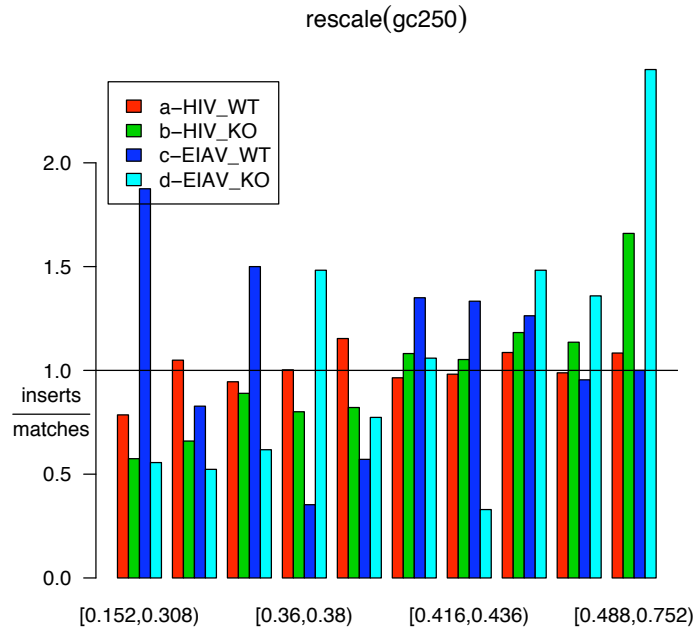


rescale(gc50)

```
            coef      se      z         p
a-HIV_WT  0.160  0.0713  2.250  2.46e-02
b-HIV_KO  0.499  0.1080  4.630  3.62e-06
c-EIAV_WT 0.102  0.2860  0.358  7.20e-01
d-EIAV_KO 0.992  0.2690  3.690  2.27e-04
```
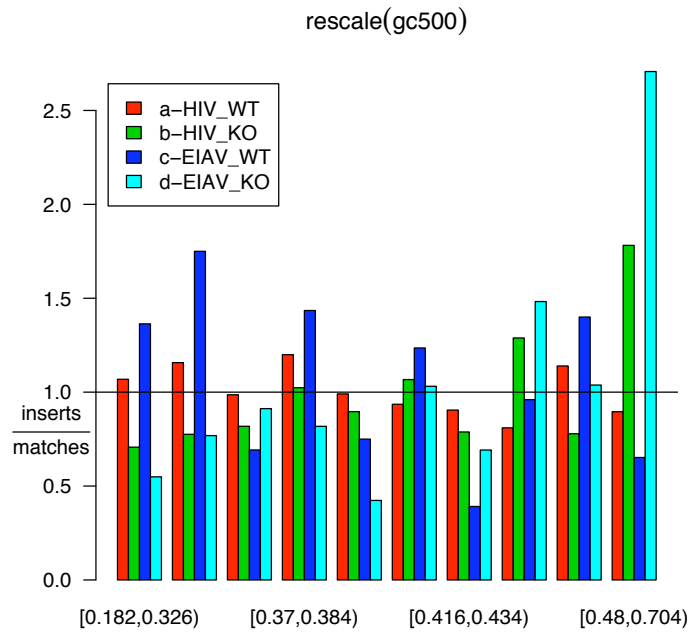
rescale(gc100)

```
            coef      se       z         p
a-HIV_WT  0.0583  0.0704  0.828  4.08e-01
b-HIV_KO  0.5590  0.1090  5.120  3.00e-07
c-EIAV_WT 0.3270  0.2770  1.180  2.39e-01
d-EIAV_KO 0.6680  0.2550  2.620  8.77e-03
```
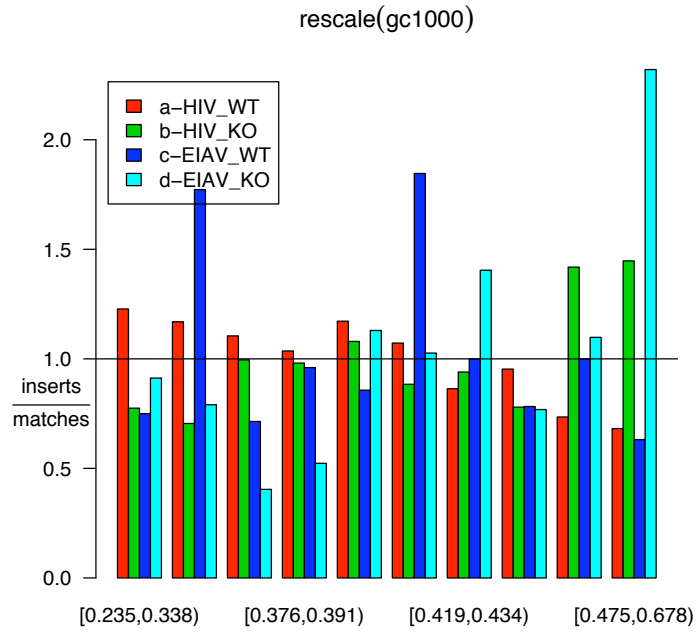


rescale(gc250)

```
                coef      se      z       p
a-HIV_WT   -0.123   0.0701  -1.75  0.08050
b-HIV_KO    0.308   0.1070   2.87  0.00408
c-EIAV_WT  -0.167   0.2740  -0.61  0.54200
d-EIAV_KO   0.660   0.2530   2.61  0.00900
```
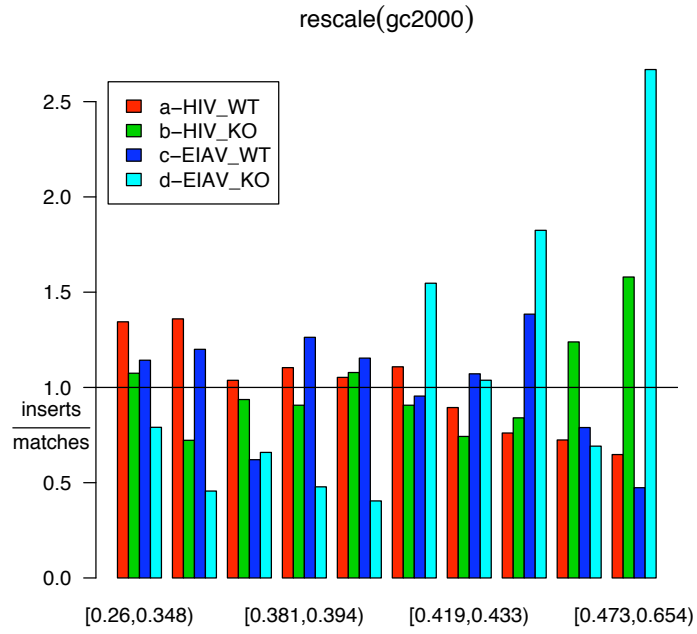
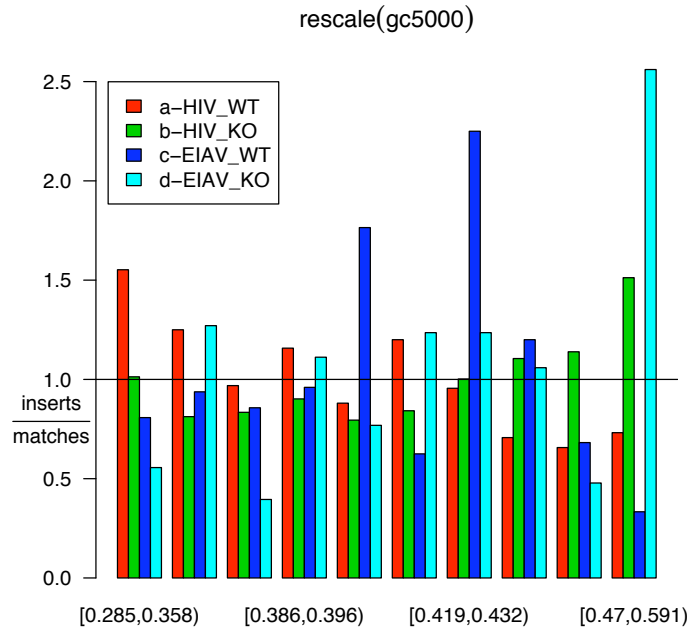rescale(gc500)

```
            coef     se        z        p
a-HIV_WT  -0.293  0.0711  -4.1200  3.79e-05
b-HIV_KO   0.193  0.1050   1.8300  6.74e-02
c-EIAV_WT -0.019  0.2760  -0.0688  9.45e-01
d-EIAV_KO  0.590  0.2540   2.3300  2.00e-02
```

rescale(gc1000)

```
              coef      se       z        p
a-HIV_WT   -0.3580  0.0709  -5.050  4.52e-07
b-HIV_KO    0.1300  0.1060   1.240  2.17e-01
c-EIAV_WT  -0.0951  0.2760  -0.344  7.31e-01
d-EIAV_KO   1.0300  0.2690   3.830  1.27e-04
```
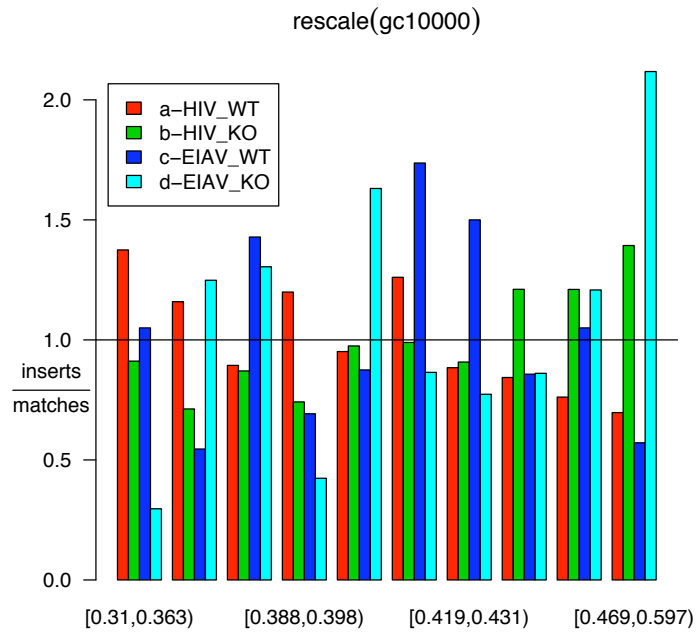
rescale(gc2000)

```
             coef      se       z        p
a-HIV_WT  -0.3180  0.0704  -4.510  6.39e-06
b-HIV_KO   0.2830  0.1060   2.670  7.60e-03
c-EIAV_WT -0.0569  0.2750  -0.207  8.36e-01
d-EIAV_KO  0.5340  0.2590   2.060  3.93e-02
```
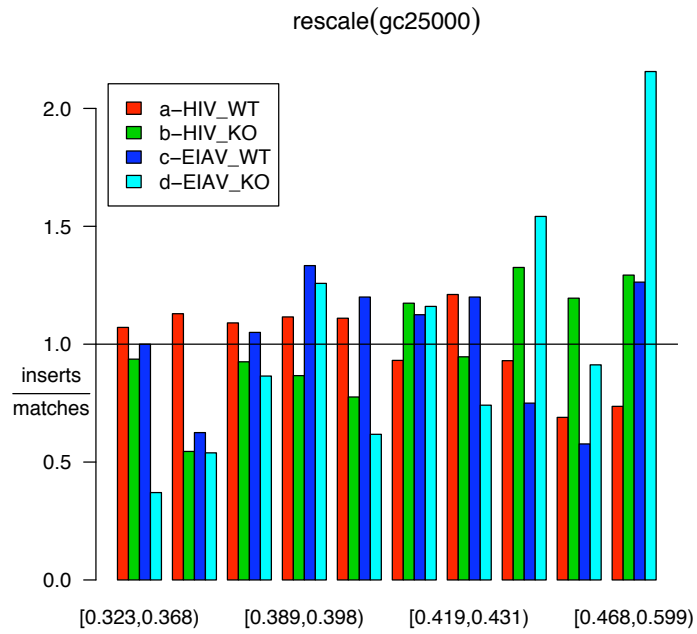


rescale(gc5000)

```
              coef     se        z         p
a-HIV_WT   -0.240  0.0705  -3.410  0.00066
b-HIV_KO    0.332  0.1070   3.110  0.00190
c-EIAV_WT   0.203  0.2720   0.746  0.45600
d-EIAV_KO   0.279  0.2500   1.120  0.26300
```
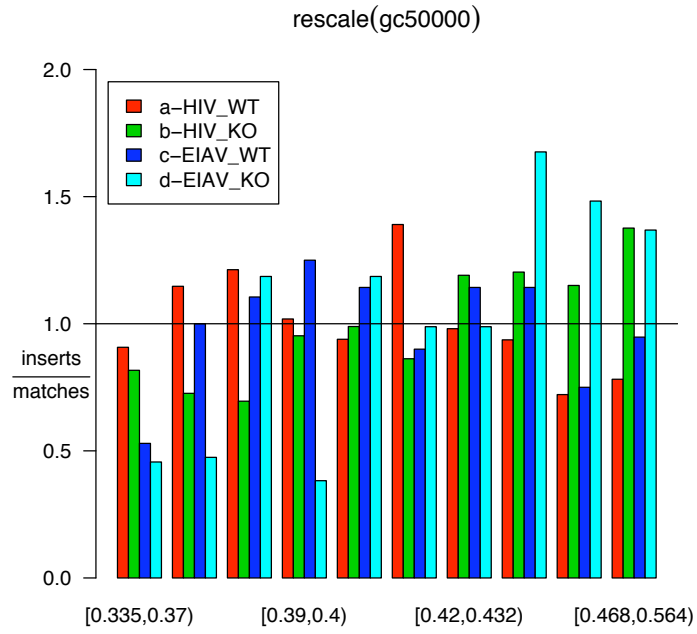
rescale(gc10000)
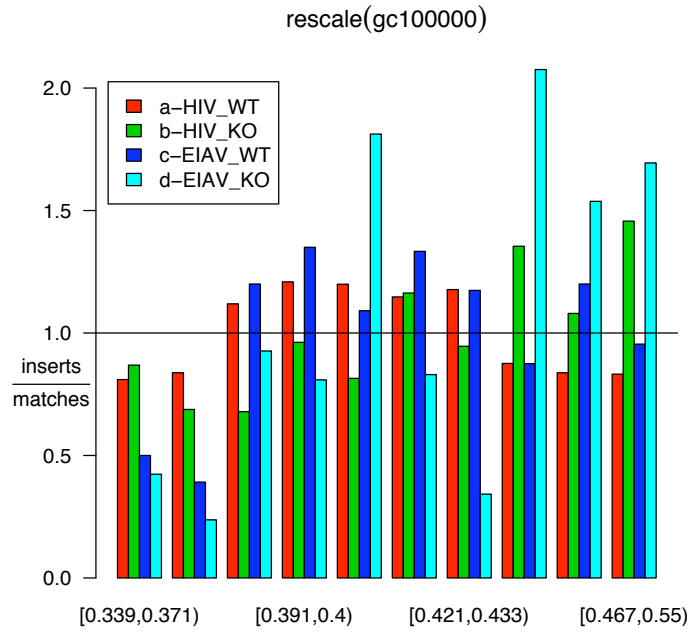
```
           coef      se       z        p
a-HIV_WT  -0.2050  0.0695  -2.950  0.003130
b-HIV_KO   0.3940  0.1070   3.670  0.000238
c-EIAV_WT -0.0931  0.2730  -0.341  0.733000
d-EIAV_KO  0.5560  0.2520   2.210  0.027300
```



rescale(gc25000)

```
              coef       se       z        p
a-HIV_WT   -0.0782   0.0694   -1.130   0.26000
b-HIV_KO    0.3400   0.1050    3.230   0.00124
c-EIAV_WT  -0.0548   0.2700   -0.203   0.83900
d-EIAV_KO   0.5990   0.2550    2.350   0.01900
```
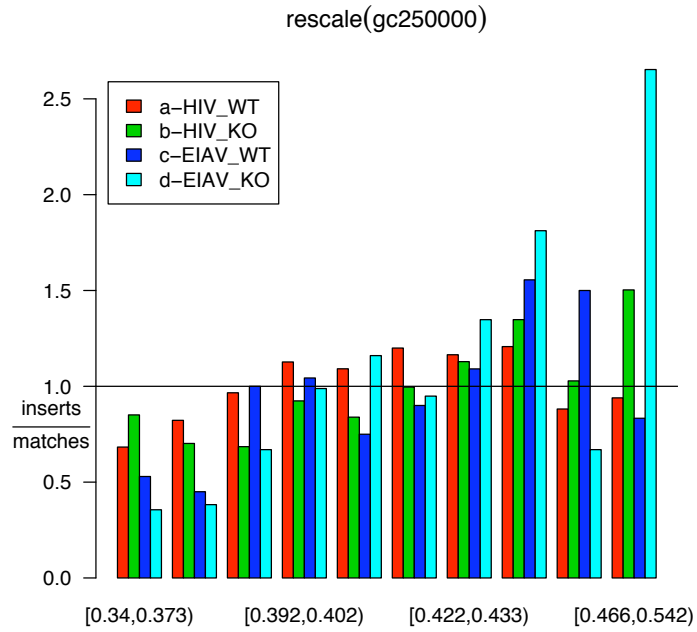


rescale(gc50000)
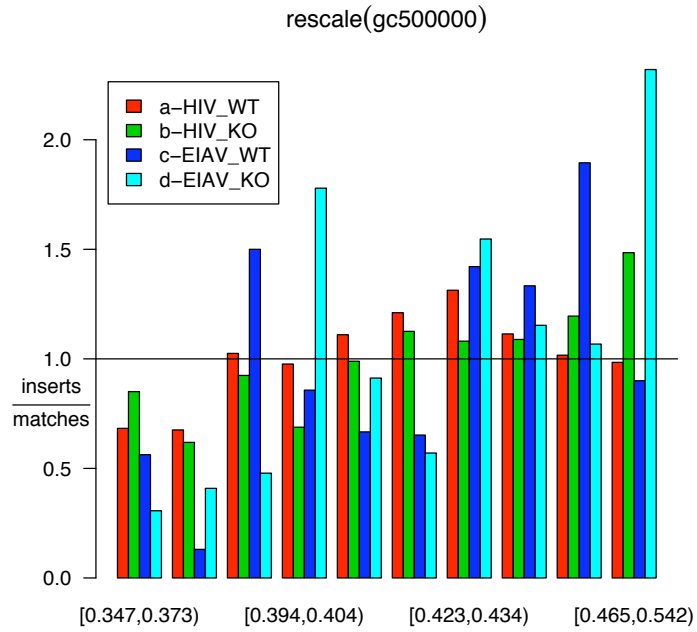
```
              coef      se       z        p
a-HIV_WT   -0.0434  0.0694  -0.625  5.32e-01
b-HIV_KO    0.4130  0.1050   3.940  8.29e-05
c-EIAV_WT   0.2030  0.2730   0.743  4.58e-01
d-EIAV_KO   0.5050  0.2570   1.960  5.00e-02
```



rescale(gc100000)
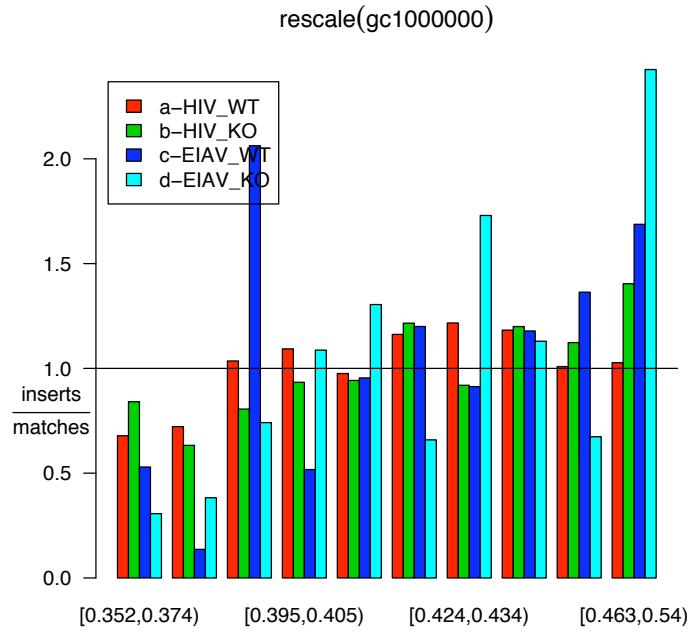
```
              coef      se      z        p
a-HIV_WT  0.142  0.0692  2.05  4.00e-02
b-HIV_KO  0.411  0.1050  3.90  9.71e-05
c-EIAV_WT 0.452  0.2920  1.55  1.22e-01
d-EIAV_KO 0.671  0.2600  2.58  9.86e-03
```

rescale(gc250000)
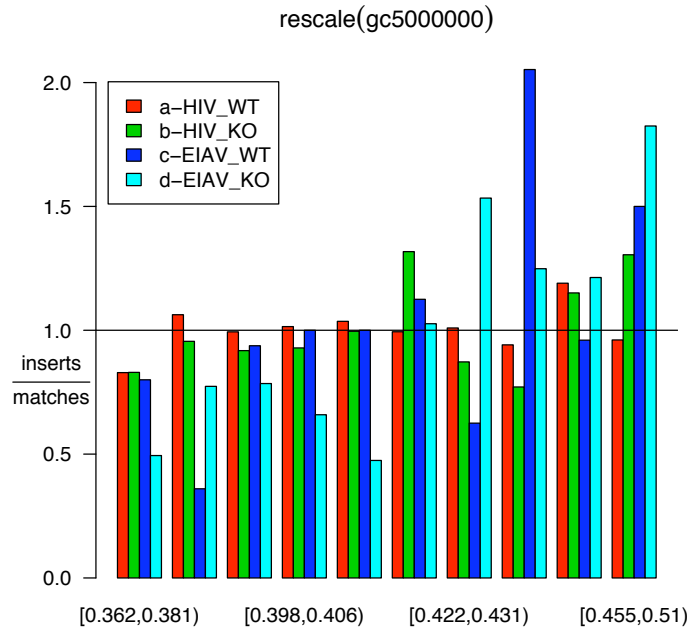
```
           coef     se    z       p
a-HIV_WT  0.246  0.0694  3.54  4.02e-04
b-HIV_KO  0.409  0.1050  3.91  9.36e-05
c-EIAV_WT 0.496  0.2940  1.69  9.16e-02
d-EIAV_KO 0.610  0.2570  2.37  1.77e-02
```



rescale(gc500000)

```
              coef      se     z        p
a-HIV_WT  0.230  0.0697  3.30  0.000971
b-HIV_KO  0.366  0.1050  3.48  0.000500
c-EIAV_WT 0.518  0.2940  1.76  0.078400
d-EIAV_KO 0.537  0.2540  2.11  0.034900
```
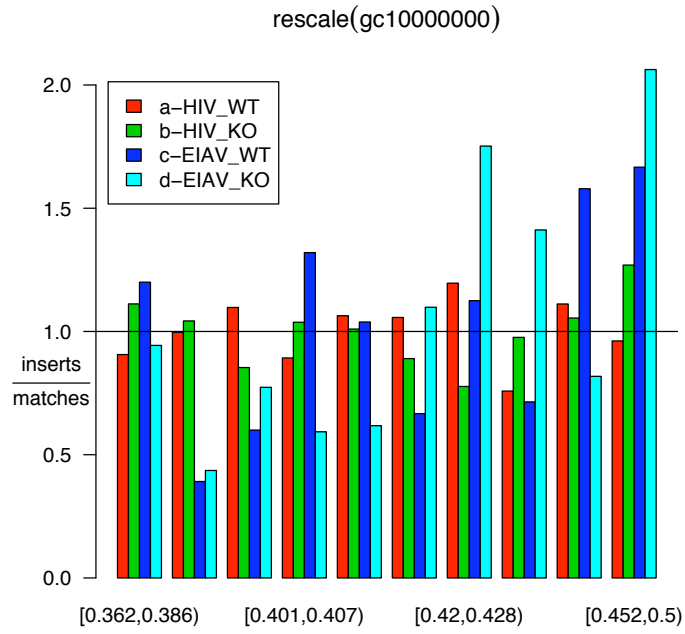
rescale(gc1000000)

```
              coef      se      z      p
a-HIV_WT  0.0339  0.0696  0.487  0.6260
b-HIV_KO  0.1590  0.1050  1.510  0.1300
c-EIAV_WT 0.4010  0.2860  1.400  0.1610
d-EIAV_KO 0.7180  0.2570  2.790  0.0052
```



rescale(gc5000000)

```
              coef      se       z         p
a-HIV_WT   0.0230  0.0696   0.331   0.74100
b-HIV_KO   0.0027  0.1040   0.026   0.97900
c-EIAV_WT  0.2030  0.2860   0.710   0.47800
d-EIAV_KO  0.7440  0.2600   2.870   0.00415
```
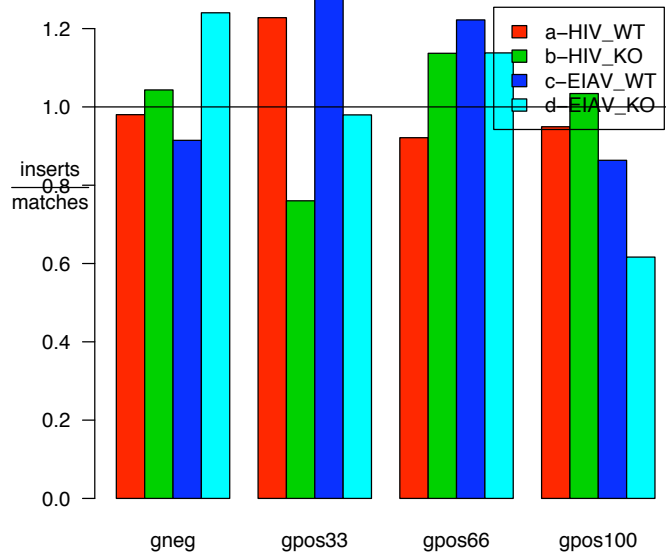
rescale(gc10000000)

# 8 Cytobands

Here we study the association of cytoband with insertion intensity. The data are obtained from
`http://genome.ucsc.edu/goldenPath/hg17/database/cytoBand.txt.gz`.



A formal test of significance attains a p-value of 0.15039. Here is the table of coefficients of the log ratio of intensities for true insertion sites versus control insertion sites (comparing each category of Giemsa staining to 'gneg') along with their standard errors, z statistics, and p-values:

```
                     coef      se       z      p
cyto.typegpos100 -0.05520  0.0664  -0.8310  0.406
cyto.typegpos33   0.06460  0.0792   0.8160  0.415
cyto.typegpos66  -0.00776  0.0915  -0.0848  0.932
```

# References

[1] Yvonne M.M. Bishop, Stephen E. Fienberg, and Paul W. Holland. *Discrete multivariate analyses: Theory and practice* (MIT Press, 1975).

[2] P. McCullagh and John A. Nelder. *Generalized linear models.* (Chapman & Hall ltd, 1999).

[3] Xiaolin Wu, Yuan Li, Bruce Crise, Shawn M. Burgess "Transcription Start Regions in the Human Genome Are Favored Targets for MLV Integration," *Science,* **300**(5626), (June 2003): 1749-1751.