# Association of Genomic Features with Integration in Stably Expressed or Inducible Cell Lines

Charles C. Berry

August 31, 2004

## Contents

# 1  Introduction

In this document, I examine the association of integration siting in cells selected as **stably expressed** (labelled 'IBB' hereafter) or **inducible** (labelled 'ID') with various genomic features.

The numbers are shown below:

```
exp.group
IBB  ID
447 388
```

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

Are there chromosomes that are particularly favored for integration by one group over the other? This was tested for statistical significance. The test performed used the likelihood ratio statistic for the logistic regression model (reviewed in [2]) as implemented by the `glm` function of R using the `binomial` family. The null hypothesis tested is the ratio of true integration events in the two groups is constant across all chromosomes. This test attains a p-value of 0.17674.

# 2 Preference for Genes

## 2.1 Acembly Genes

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



Is there is a difference in the tendency for insertions to occur in genes? A formal test of significance yields a p-value of 0.053439. In the following plot we show the relative frequency of insertions in exons according to the 'Acembly'

anotation 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 ID sites versus IBB sites along with their standard errors, z statistics, and p-values:

```
              coef    se      z       p
(Intercept)  0.166 0.174   0.953 0.3410
in.gene     -0.426 0.193  -2.210 0.0270
in.exon      0.391 0.211   1.860 0.0632
```

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 the both introns and intergenic regions, so the impression given by the table may differ from that for the barplot.

## 2.2   refGenes

Here we examine the relative preference that insertions of the two types have for genes. In the following plot we show the relative frequency of insertions in genes according to the 'refGene' anotation.

Is there is a tendency for insertions to occur in genes? A formal test of significance yields a p-value of 0.86057.

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

Here is the table of coefficients of the log ratio of intensities for ID sites versus IBB sites along with their standard errors, z statistics, and p-values:

```
                coef    se      z      p
(Intercept) -0.1570 0.112 -1.4000 0.162
in.gene      0.0137 0.144  0.0947 0.925
in.exon      0.2180 0.396  0.5500 0.583
```

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 the both introns and intergenic regions, so the impression given by the table may differ from that for the barplot.

## 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' anotation.

Is there is a tendency for insertions to occur in genes? A formal test of significance yields a p-value of 0.091842.

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

Here is the table of coefficients of the log ratio of intensities for ID sites versus IBB sites along with their standard errors, z statistics, and p-values:

```
              coef    se       z      p
(Intercept)   0.0741 0.146   0.509 0.611
in.gene      -0.2890 0.166  -1.740 0.082
in.exon       0.3110 0.444   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. Note that in the barplot above the 'Not in Exon' bars include the both introns and intergenic regions, so the impression given by the table may differ from that for the barplot.

data set.

## 2.4 uniGenes

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 'uniGene' anotation.

Is there is a tendency for insertions to occur in genes? A formal test of significance yields a p-value of 0.46991.

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

Here is the table of coefficients of the log ratio of intensities for ID sites versus IBB sites along with their standard errors, z statistics, and p-values:

```
              coef    se      z      p
(Intercept) -0.0606 0.132 -0.460 0.645
in.gene     -0.1210 0.157 -0.769 0.442
in.exon      0.0863 0.267  0.324 0.746
```

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 the both introns and intergenic regions, so the impression given by the table may differ from that for the barplot.

# 3 CpG Island Neighborhoods

Here we study the effect of being in the neighborhood of CpG Islands. Following Wu et al [4], who found that the neighborhoods within $\pm1$kb 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 $\pm1$kb of a CpG island:

A formal test of significance comparing the difference attains a p-value of 0.55736.

## 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 0.09008.

## 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 0.17307.

## 3.4   25 kilobase neighborhoods

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

A formal test of significance comparing the difference attains a p-value of 0.39436.

## 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 0.28185.

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

In this section the association with gene density is examined. The 'genes' that are counted are the Ensembl genes. In addition, we study various functions of the EST counts for the Ensembl genes using data described in Versteeg et al [3] and CpG Island density. Based on preliminary observations, it was decided to determine the density of ESTs found in a region in the following ways:

**count.exprs**  Count only one EST per gene and divide by number of bases

**exprs**  Count up to 200 ESTs per gene and divide by number of bases

**big.exprs**  Counting only the ESTs in excess of two hundred per gene 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 often 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, then 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 quadratic polynomial to the gene density values.

**dens.25k  – p–value = 0.34438**



In the barplot that follows we examine the association of insertion sites with expression density in a 25 kilobase window surrounding each locus. First, we count just one EST per gene.

16

**count.exprs.25k – p–value = 0.40539**

Now we count up to 200 ESTs per gene:

**exprs.25k – p–value = 0.13277**

And here counting starts only after 200 ESTs per gene

**big.exprs.25k – p–value = 0.64815**

Here the effect of density of CpG islands is studied:

**cpg.dens.25k – p–value = 0.073691**

## 4.2   50 kiloBase Window

First, we see gene density:

**dens.50k – p–value = 0.045084**

Here are the results for EST density. First, we count just one EST per gene.

**count.exprs.50k – p–value = 0.13572**

Now we count up to 200 ESTs per gene:

**exprs.50k  – p–value = 0.018256**

And here counting starts only after 200 ESTs per gene

**big.exprs.50k – p–value = 0.30184**

Here the effect of density of CpG islands is studied:

**cpg.dens.50k – p–value = 0.043417**



## 4.3   100 kiloBase Window

First, we see gene density:

**dens.100k – p–value = 0.18165**

Here are the results for EST density. First, we count just one EST per gene.

**count.exprs.100k – p–value = 0.40178**

Now we count up to 200 ESTs per gene:

**exprs.100k – p–value = 0.10123**

And here counting starts only after 200 ESTs per gene

**big.exprs.100k – p–value = 0.20328**

Here the effect of density of CpG islands is studied:

## 4.4   250 kiloBase Window

First, we see gene density:

**dens.250k – p–value = 0.33382**

Here are the results for EST density. First, we count just one EST per gene.

**count.exprs.250k – p–value = 0.060095**

Now we count up to 200 ESTs per gene:

**exprs.250k  – p–value = 0.5713**

And here counting starts only after 200 ESTs per gene

**big.exprs.250k – p–value = 0.38762**

Here the effect of density of CpG islands is studied:

## 4.5   500 kiloBase Window

First, we see gene density:

**dens.500k – p–value = 0.10433**

Here are the results for EST density. First, we count just one EST per gene.

**count.exprs.500k – p–value = 0.039490**

Now we count up to 200 ESTs per gene:

**exprs.500k – p–value = 0.24795**

And here counting starts only after 200 ESTs per gene

**big.exprs.500k – p–value = 0.21487**

Here the effect of density of CpG islands is studied:

**cpg.dens.500k – p–value = 0.00032722**

## 4.6 1 megaBase Window

First, we see gene density:

**dens.1M – p–value = 0.097193**

Here are the results for EST density. First, we count just one EST per gene.

**count.exprs.1M – p–value = 0.010932**

Now we count up to 200 ESTs per gene:

**exprs.1M  – p–value = 0.21409**

And here counting starts only after 200 ESTs per gene

**big.exprs.1M – p–value = 0.47466**

Here the effect of density of CpG islands is studied:

**cpg.dens.1M – p–value = 0.0018961**

## 4.7  2 megaBase Window

First, we see gene density:

**dens.2M  – p–value = 0.0385**

Here are the results for EST density.  First, we count just one EST per gene.

**count.exprs.2M – p–value = 0.0091795**

Now we count up to 200 ESTs per gene:

**exprs.2M – p-value = 0.21620**

And here counting starts only after 200 ESTs per gene

**big.exprs.2M – p–value = 0.81479**

Here the effect of density of CpG islands is studied:

**cpg.dens.2M – p–value = 0.0057697**

## 4.8   4 megaBase Window

First, we see gene density:

**dens.4M – p–value = 0.35540**

Here are the results for EST density. First, we count just one EST per gene.

**count.exprs.4M – p–value = 0.087686**

Now we count up to 200 ESTs per gene:

**exprs.4M – p–value = 0.064126**

And here counting starts only after 200 ESTs per gene

**big.exprs.4M – p–value = 0.64414**

Here the effect of density of CpG islands is studied:

**cpg.dens.4M – p–value = 0.0019046**

## 4.9   4 megaBase Window

First, we see gene density:

**dens.8M – p-value = 0.58613**

Here are the results for EST density. First, we count just one EST per gene.

**count.exprs.8M – p–value = 0.40197**

Now we count up to 200 ESTs per gene:

**exprs.8M – p–value = 0.30922**

And here counting starts only after 200 ESTs per gene

## 4.10   16 megaBase Window

First, we see gene density:

**dens.16M – p-value = 0.33148**

Here are the results for EST density. First, we count just one EST per gene.

**count.exprs.16M  – p–value = 0.19284**



Now we count up to 200 ESTs per gene:

And here counting starts only after 200 ESTs per gene

## 4.11 32 megaBase Window

First, we see gene density:

**dens.32M – p–value = 0.059183**

Here are the results for EST density. First, we count just one EST per gene.

**count.exprs.32M  – p–value = 0.011657**

Now we count up to 200 ESTs per gene:

**exprs.32M – p–value = 0.15809**

And here counting starts only after 200 ESTs per gene

**big.exprs.32M – p–value = 0.30093**

# 5 Juxtaposition with Gene Start and End Positions

## 5.1 Acembly 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 Acembly gene.

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

**acembly non–gene width  – p–value = 0.00014776**

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.

**acembly boundary.dist – p−value = 0.33025**



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.

69

**acembly start.dist – p–value = 0.060898**

## 5.2   RefSeq Annotations

**refSeq non−gene width  – p−value = 1.1641e−06**

**refSeq boundary.dist – p–value = 0.7536**

refSeq start.dist  – p−value = 0.25313

## 5.3   genScan Annotations

**genScan non−gene width  – p−value = 0.00068707**

**genScan boundary.dist  – p–value = 0.9533**

genScan start.dist – p–value = 0.34622

## 5.4   uniGene Annotations

**uniGene non–gene width  – p–value = 0.00068707**

uniGene boundary.dist  – p–value = 0.9533

**uniGene start.dist  – p–value = 0.34622**

## 6  GC content

Here we study the effect of GC content on insertion. The GC content is taken from the Human Genome Draft at GoldenPath from the table
`http://genome.ucsc.edu/goldenPath/14nov2002/database/gcPercent.txt.gz`.

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

**gcpct – p–value = 0.526**



# 7 Cytobands

Here we study the association of cytoband with insertion intensity. The data are obtained from
`http://genome.ucsc.edu/goldenPath/14nov2002/database/cytoBand.txt.gz`.

A formal test of significance attains a p-value of 0.41588.

# 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] Rogier. Versteeg, Barbera. D. C. van Schaik., Marinus. F. van Batenburg., Marco. Roos, Ramin. Monajemi, Huib. Caron, Harmen. J. Bussemaker, and Antoine. H. C. van Kampen. The human transcriptome map reveals extremes in gene density, intron length, GC content, and repeat pattern for domains of highly and weakly expressed genes. *Genome Res*, 13(9):1998–2004, Sep 2003.

[4] Xiaolin. Wu, Yuan. Li, Bruce. Crise, and Shawn. M. Burgess. Transcription start regions in the human genome are favored targets for MLV integration. *Science*, 300(5626):1749–1751, Jun 2003.