

# A Modification Motif Analysis Demo

Pat Marks

April 29, 2012

## 1 Introduction

This document demonstrates a typical analysis of kinetic modification detection data from bacterial genome. The raw data used here is generated by sequencing native E.Coli K12 on the PacBio RS, and performing secondary analysis using the RS.Modification.Detection workflow. Some words about RM dogma – decoding the RM system is the goal of this analysis.

- Load and explore the modification calls from the RS.Modification.Detection workflow
- Find the sequence motifs that are common to most of the modified genome hits
- Find instances of those motifs in the genome sequence
- Determine the what fraction of the motifs instances were detected as methylated

### 1.1 R Setup

The analysis workflow demonstrated here on R 2.14 and R 2.15 on Linux and Windows. From CRAN we require `ggplot2` and `plyr`. From Bioconductor we require `Biostrings` and `cosmo`. See the appendix for details of download and installation details.

Please prepare your R environment as follows. `scripts.R` is included with this tutorial.

```
> library(ggplot2)
> library(plyr)
> library(Biostrings)
> library(cosmo)

Welcom to cosmo version 1.18.0

cosmo is free for research purposes only. For more details, type
license.cosmo(). Type citation('cosmo') for details on how to cite
cosmo in publications.

> source("scripts.R")
```

## 2 Loading modification data

We start by loading the raw modifications calls from that are produced by SMRTPortal in `modifications.gff.gz`. This can be downloaded from the job results page in SMRTPortal, or accessed directly from the SMRTportal job folder on the file server. `modifications.gff.gz` is compliant with GFFv3 specification ([www.sequenceontology.org]). See Table ?? for a description of the columns in the GFF file.

The included R script contains a GFF file reader that extracts some extra attribute columns used by the modification detection tool. Take a look at the data contained in the GFF file. The context field contains a 41 base context centered around the detected modification – pull out the center base (position 21). Summarize the `coverage` column of the GFF to see how many reads contribute to each modification call. Confident m6A detection generally requires coverage > 20 per strand.

Column	Description
seqid	Reference tag (e.g. ref00001)
source	Name of tool – ‘kinModCall’
type	Modification type – currently we use a generic tag ”modified_base”
start	Location of modification
end	Location of modification plus one
score	Phred transformed p-value of detection
strand	Sample strand containing modification
phase	Not applicable
attributes	Fields below are packed in the GFF attributes column
IPDRatio	Ratio between mean IPD of observed data to IPD of unmodified DNA
context	Reference sequence -20bp to +20bp around start, converted to current strand
coverage	Number of valid IPD observations at this site

Table 1: Contents of modifications.gff.gz file

```
> gff <- "/mnt/secondary/Smrtanalysis/userdata/jobs/040/040041/data/modifications.gff.gz"
> hits <- readModificationsGff(gff)

Reading /mnt/secondary/Smrtanalysis/userdata/jobs/040/040041/data/modifications.gff.gz

> hits$CognateBase <- substr(hits$context, 21, 21)
> head(hits)

  seqname      source      feature start end score strand frame coverage
1 ref000001 kinModCall modified_base 271 271   26     -    . 34
2 ref000001 kinModCall modified_base 423 423   21     -    . 36
3 ref000001 kinModCall modified_base 621 621   81     -    . 37
4 ref000001 kinModCall modified_base 653 653   21     -    . 35
5 ref000001 kinModCall modified_base 728 728   65     -    . 40
6 ref000001 kinModCall modified_base 738 738   30     -    . 39

           context IPDRatio CognateBase
1 TCAGGTGGGGCTTTTCTGTGTTCTGTACCGTCAGC 3.43      G
2 ACGGTGGCACCTGCCCTGCCTGGCATTGCTTCCAGAAT 2.04      C
3 TTTATTTGGCAAATTCCCTGATCGACGAAAGTTCAATTG 5.18      A
4 GCCCAACAAAATAATGCCATGCAGGACATGTTTATTG 1.83      T
5 ATACGCCGCCATAATGGCGATCGACATTTCTGCCACGG 2.83      A
6 CGCGCTTCTAACAGCCGCCATAATGGCGATCGACATTG 1.95      C

> summary(hits$coverage)

  Min. 1st Qu. Median Mean 3rd Qu. Max.
  5.00   41.00  47.00  47.86  54.00 223.00
```

We now make some plots from the GFF data to assess the quality and type of the modification calls. Figure 1 shows a histogram of the scores of the GFF entries, coloured by the cognate base. This will give you a sense of how strong your signals are and whether the strongest signals are enriched on any base. For our E.Coli test genome the predominant modification is 6-methyl adenosine, so most of the significant modification detections are at A positions.

The histogram in Figure 1 indicates that the interesting A bases have a score cutoff of roughly 45. We select these hits, then sort in decreasing order of score, so we consider the strongest signal first.

```
> goodHits <- subset(hits, score > 45)
> goodHits <- goodHits[order(goodHits$score, decreasing = T), ]
> workHits <- goodHits
```

```
> p <- qplot(score, colour = CognateBase, geom = "freqpoly", data = hits,  
+           binwidth = 5)  
> show(p)
```

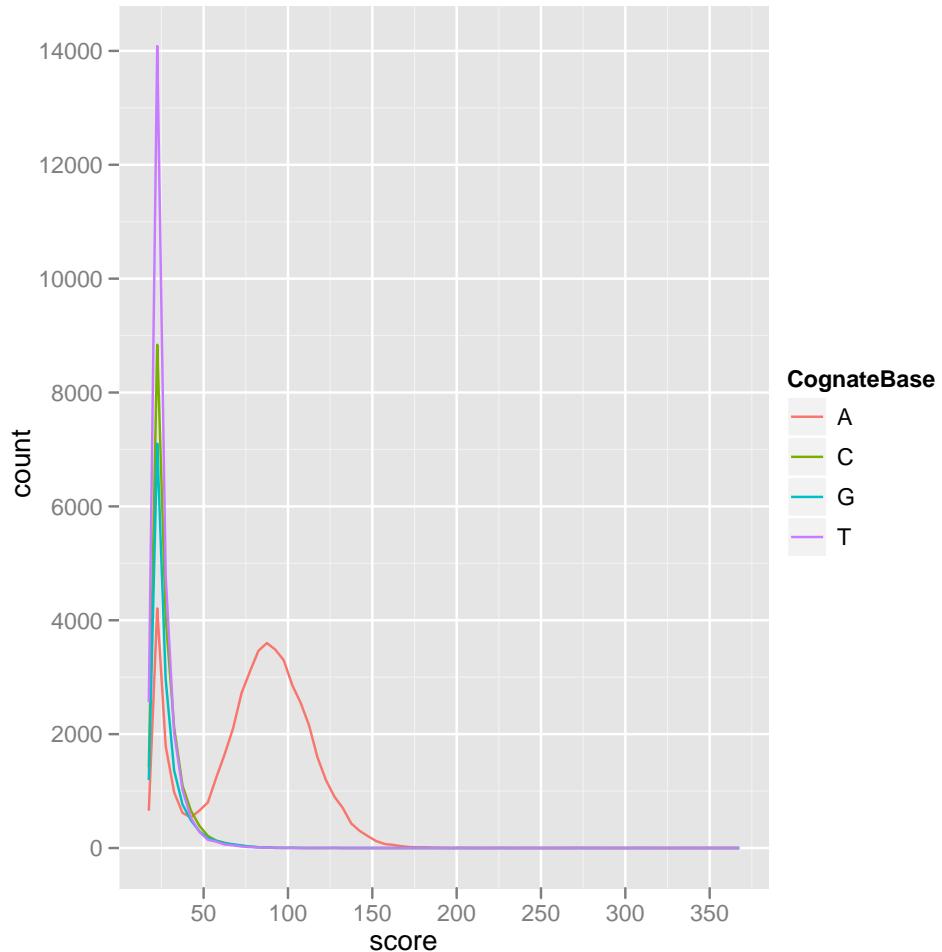


Figure 1: Modification Scores by cognate base

```
> p <- qplot(coverage, score, colour = CognateBase, alpha = I(0.3),  
+             data = hits)  
> show(p)
```

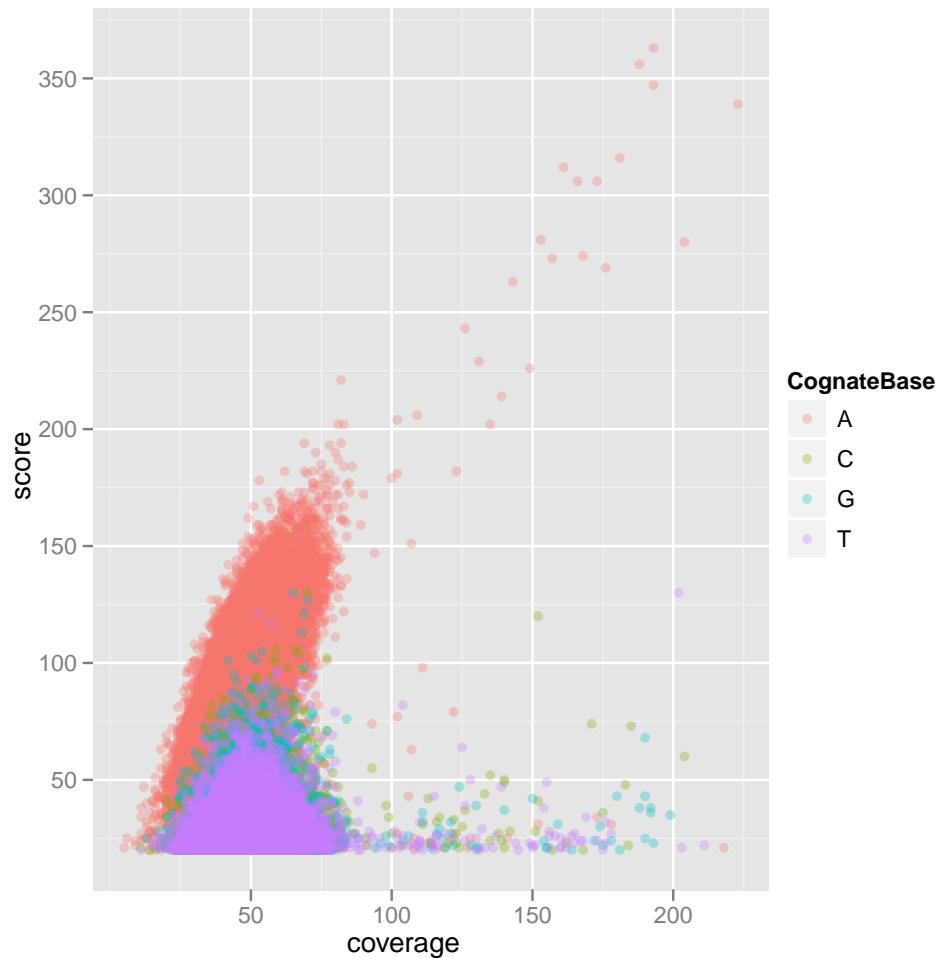


Figure 2: Score vs. Coverage

### 3 Motif Finding

We use the motif finding package `cosmo` (<http://www.bioconductor.org/packages/release/bioc/html/cosmo.html>) to search for conserved motifs in the DNA sequences surrounding our modification detections. First, we convert the context snippets into a `cosmo` compatible format. Next, we select a small number of hits to pass to `cosmo` (it will not return if you provide too many at once). We feed the contexts to `cosmo`, then look at the resulting logo plot (Figure 3. The 'GATC' motif is strongly detected. `cosmo` searches for motifs anywhere within the supplied sequences. In our case the motif is always appears at the same position within sequence windows we supplied. The motifs hits detected by `cosmo` should therefore have the same starting position. We check for this by examining `table(cosmoOut@motifs@pos)`. Clearly the mode is at position 15. The A in the GATC motif occurs at position 7 in the `cosmo` motif, which corresponds to  $15 + 7 - 1 = 21$ , the center, cognate base of the sequences we passed to `cosmo`

```
> sequenceList <- lapply(workHits$context[1:200], function(x) list(seq = as.character(x),
+     desc = ""))
> cosmoOut <- cosmo(seqs = sequenceList, minW = 4, maxW = 10, revComp = F,
+     silent = T)
> table(cosmoOut@motifs@pos)

 2   4   6   8   9   10  15  24  25  31
 2   1   3   1   1    2 185   1   3   1

> startPos <- 15 + 7 - 1
> startPos

[1] 21
```

We view the logo plot of the motif found (Figure ??)

We can check which hits match this context using the `labelContexts` function and see how many GATC hits we have. The `labelContexts` takes a charactervector of context strings, and a vector of motif strings and their associated methylated positions. We make a rough approximation of the number of GATC sites we would expect on both stands of a random genome of length  $|G|$  as  $2|G|/4^n$  where  $n$  is the number of bases in the motif. We see that most GATC are methylated. We will perform a more accurate analysis in the next section. We can now remove the 'GATC' hits from working list. We still have many of unassigned GFF hits, so we can run `cosmo` on the reduced list of hits.

```
> motif <- "GATC"
> position <- 2
> motifLabels <- labelContexts(workHits$context, motif, position)
> table(motifLabels)

motifLabels
  GATC  None
37728 4192

> genomeSize <- max(hits$end)
> expectedHits <- 2 * genomeSize/(4^4)
> expectedHits

[1] 36245.39

> workHits <- workHits[motifLabels != motif, ]
> nrow(workHits)

[1] 4192
```

We now run `cosmo` on the remaining:

```
> plot(cosmoOut)
```

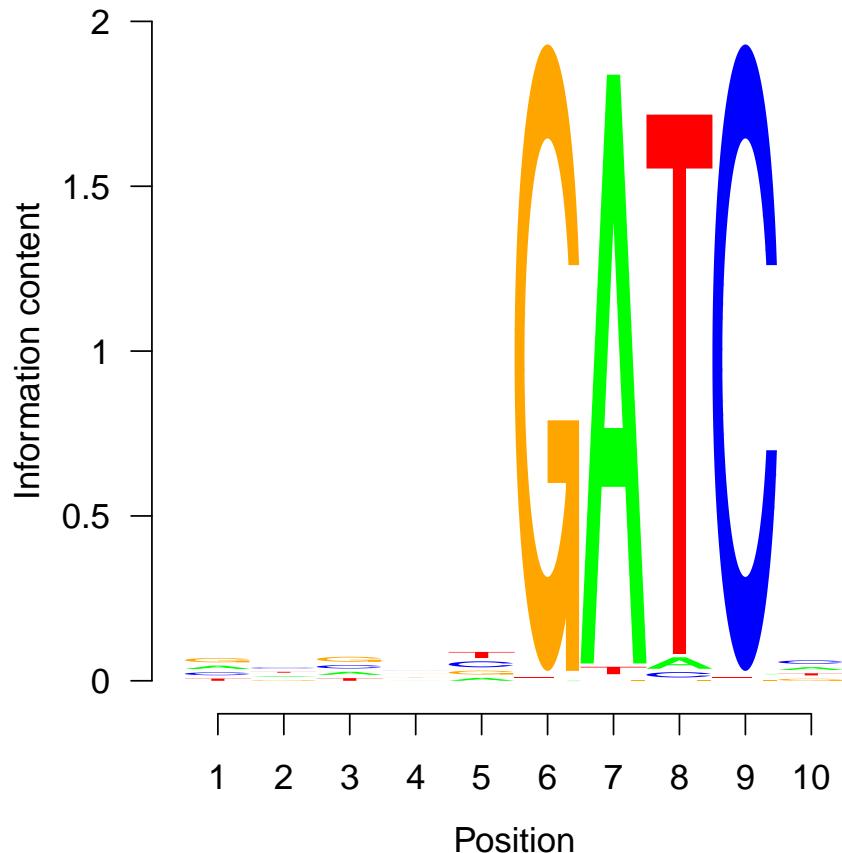


Figure 3: Logo plot - First `cosmo` iteration

```
> plot(cosmoOut2)
```

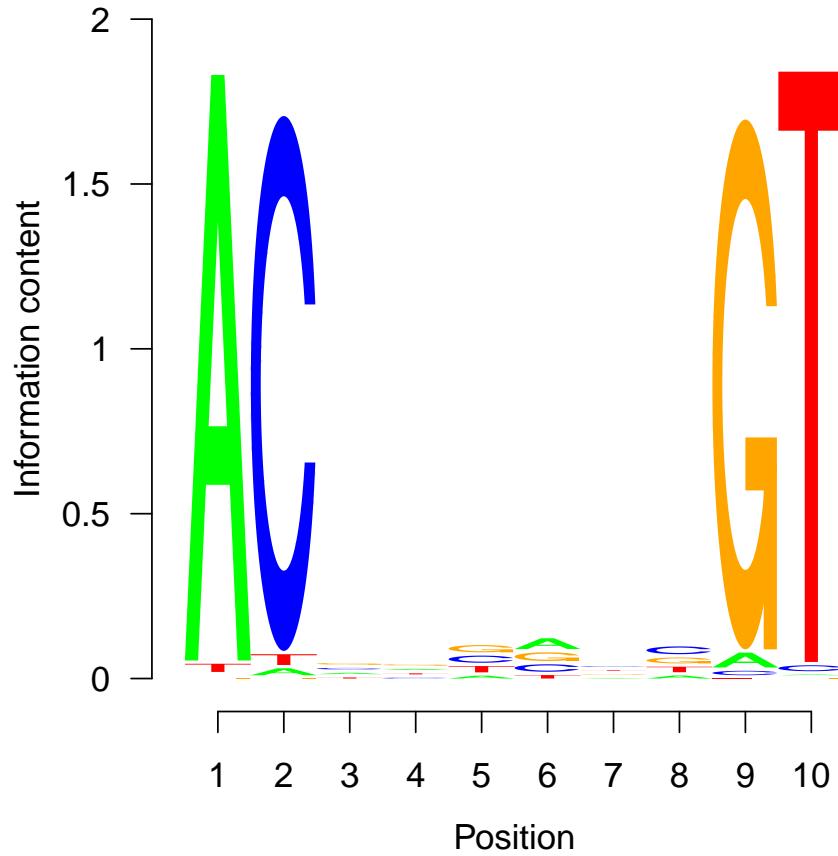


Figure 4: Logo plot: Second cosmo iteration

```
> sequenceList <- lapply(workHits$context[1:200], function(x) list(seq = as.character(x),
+     desc = ""))
> cosmoOut2 <- cosmo(seqs = sequenceList, minW = 4, maxW = 10,
+     revComp = F, silent = T)
> motif <- "ACNNNNNNGT"
> position <- 1
> motifLabels <- labelContexts(workHits$context, motif, position)
> table(motifLabels)

motifLabels
ACNNNNNNGT      None
 1183      3009

> genomeSize <- max(hits$end)
> expectedHits <- 2 * genomeSize/(4^4)
> expectedHits

[1] 36245.39
```

Figure 4 shows the logo plot of the second cosmo run. It's unclear what's happening now. We have found a highly confident 4 base motif but we observe far less than the expected number. In fact our E.Coli sample

contains a pair of methyltransferases that methylate reverse-complementary motifs that share a common 4 base sub-motif. These two motifs occur in equal numbers. `cosmo` was not designed to handle situation like this, so we must continue on by hand, inspecting the hits manually. Hand inspection of the hits shows that 'GATC', 'GCACNNNNNGTT', 'AACNNNNNGTGC' are the fully methylated motifs in our sample.

## 4 Genome Annotation

We can load the genome sequence and annotate it instances of each motif, to determine the genome-wide methylation fraction of our motifs. The `genomeAnnotation` function returns a `data.frame` containing one row for each match in the supplied genome of each motif supplied. Genome positions can match multiple motifs, which gives multiples row at the same genome position

```
> seq_path <- "/mnt/secondary/Smrtanalysis/userdata/references/ecoli/sequence/ecoli.fasta"
> dna_seq <- read.DNAStringSet(seq_path)
> motifs = c("GATC", "GCACNNNNNGTT", "AACNNNNNGTGC")
> positions = c(2, 3, 2)
> genomeAnnotations <- genomeAnnotation(dna_seq, motifs, positions)
> head(genomeAnnotations)

  strand start motif onTarget seqid
1      +    620  GATC      On     1
2      +    727  GATC      On     1
3      +    782  GATC      On     1
4      +    881  GATC      On     1
5      +   1168  GATC      On     1
6      +   1570  GATC      On     1

> table(genomeAnnotations$motif)

AACNNNNNGTGC      GATC GCACNNNNNGTT
      595          38240         595
```

We merge the `genomeAnnotation` output with our `GFF` `data.frame` to count the number of motif instances that exist in the genome and the number that were detected as methylated. We merge the `genomeAnnotation` and `goodHits` tables by genome position and strand, with `all=T` to include `GFF` hits that are not annotated with a motif, and genome motif instances that do not have a `GFF` hit. We adjust the merged `data.frame` to indicate these cases.

```
> goodHits$seqid <- as.integer(substr(goodHits$seqname, 4, 11))
> mm <- merge(goodHits, genomeAnnotations, all = T)
> mm$motif[is.na(mm$motif)] <- "NoMotif"
> mm$feature[is.na(mm$feature)] <- "not_detected"
> table(mm$feature, mm$motif)

  AACNNNNNGTGC  GATC GCACNNNNNGTT NoMotif
modified_base      576 37728        584    3032
not_detected       19   512         11      0
```

In our `goodHits` table we have 3032 'modified\_base' calls that do not occur at an annotated genome position. For genome positions matching the GATC motif, 37728 were present detected ('modified\_base') and 512 were not ('not\_detected'). We can assess whether we are likely to have missed any methylated motifs by comparing the score distributions for the positions matching a motif to the 'NoMotif' set, as in Figure 5.

## A R Package Installation

`R` can be downloaded for Linux, Mac and Windows from <http://r-project.org>. We also recommend the graphical front-end `RStudio`, available from <rstudio.org>. The following `R` packages need to be installed to run through this demo. `ggplot2`, `plyr` are available through the built-in `R` package manager. Instructions for installing Bioconductor packages is available here: <www.bioconductor.org/install>. `cosmo` and `Biostrings` are required from Bioconductor

```

> p <- qplot(score, ..density.., colour = motif, geom = "freqpoly",
+   data = subset(mm, feature == "modified_base"), binwidth = 3,
+   xlim = c(0, 200))
> show(p)

```

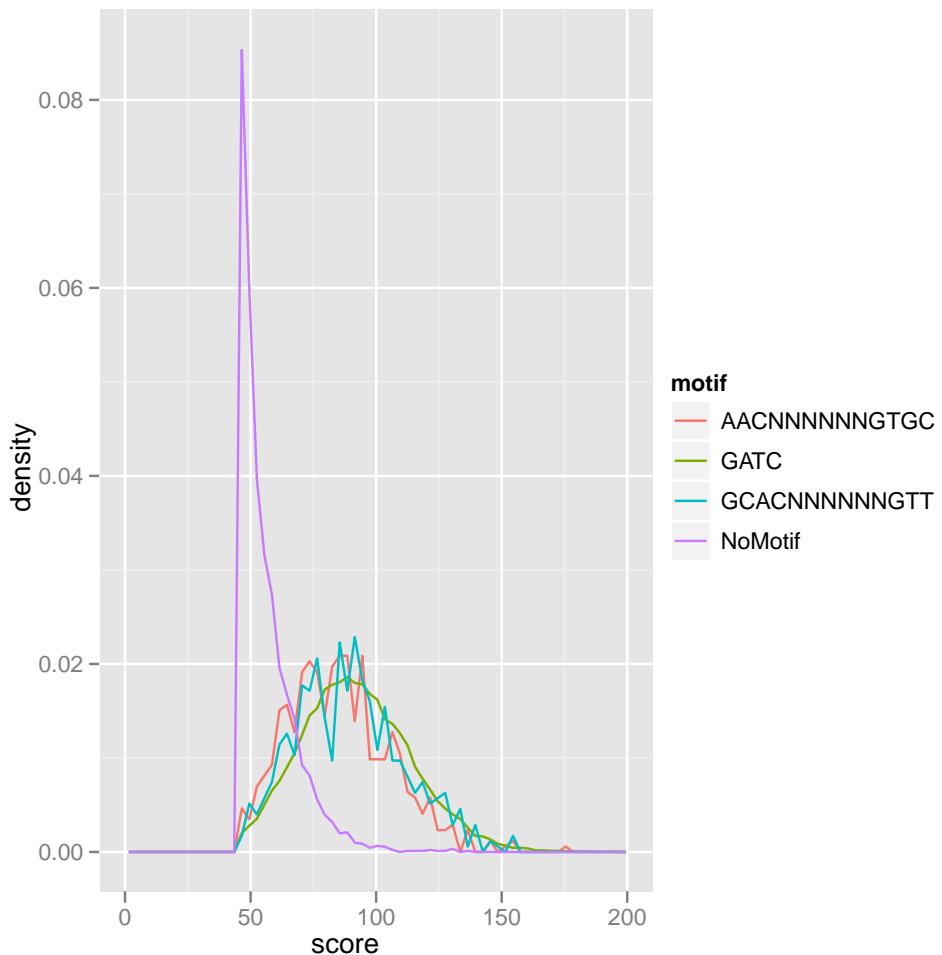


Figure 5: Score distribution by motif annotation

## A.1 Session Info

Here we give information about the version of R and installed packages used to generate this document.

```
> sessionInfo()

R version 2.13.1 Patched (2011-09-13 r57007)
Platform: x86_64-unknown-linux-gnu (64-bit)

locale:
[1] LC_CTYPE=en_US.UTF-8      LC_NUMERIC=C
[3] LC_TIME=en_US.UTF-8       LC_COLLATE=en_US.UTF-8
[5] LC_MONETARY=C            LC_MESSAGES=en_US.UTF-8
[7] LC_PAPER=en_US.UTF-8     LC_NAME=C
[9] LC_ADDRESS=C              LC_TELEPHONE=C
[11] LC_MEASUREMENT=en_US.UTF-8 LC_IDENTIFICATION=C

attached base packages:
[1] grid      stats      graphics grDevices utils      datasets methods
[8] base

other attached packages:
[1] cosmo_1.18.0    seqLogo_1.18.0   Biostrings_2.20.4 IRanges_1.10.6
[5] ggpplot2_0.8.9   proto_0.3-9.2    reshape_0.8.4    plyr_1.7.1

loaded via a namespace (and not attached):
[1] digest_0.5.1 tools_2.13.1
```

For Research Use Only. Not for use in diagnostic procedures. Copyright 2011, Pacific Biosciences of California, Inc. All rights reserved. Information in this document is subject to change without notice. Pacific Biosciences assumes no responsibility for any errors or omissions in this document. Certain notices, terms, conditions and/or use restrictions may pertain to your use of Pacific Biosciences products and/or third party products. Please refer to the applicable Pacific Biosciences Terms and Conditions of Sale and to the applicable license terms at <http://www.pacificbiosciences.com/licenses.html>.

Pacific Biosciences, the Pacific Biosciences logo, PacBio, SMRT and SMRTbell are trademarks of Pacific Biosciences in the United States and/or certain other countries. All other trademarks are the sole property of their respective owners.