

# Hw2

## Home work 2 - to be done as groups

Names: Adrian Ramon Santonja, Jonas B H Andersen, Ismael Rodriguez Palomo, Steen J. Østergaard, Luiza Czerwinska.

Group: 3

For deadlines etc, see Absalon.

### Question 1: Dicer dissected

The human DICER1 gene encodes an important ribonuclease, involved in miRNA and siRNA processing. Several mRNAs representing this gene have been mapped to the human genome (March 2006 assembly). We will look closer at one of them with the accession number AK002007.

a) What are the first five genomic nucleotides that are read by RNA polymerase II from this transcript?

**Answer:** The first 5 genomic nucleotides seen from the UCSC-genome browser is: AAAGG

This is seen on the following screenshot (Fig 1.) with the first exon starting on the right and running to the left. The sequence CCTTT is reverse complemented and gives AAAGG.

b) Look at the raw mRNA sequence of AK002007, from the database it actually comes from. What are the first five nucleotides?

The first 5 nucleotides from the GenBank sequence is gaagcaa. This is seen in the screenshot on figure 2:

c) How do you explain the discrepancy (maximum 5 lines)?

The discrepancy is hard to explain, but we have a couple of theories. Looking at the GenBank entry we can see that the sequences are found by oligo-capping. In this method a cDNA library is constructed by removal of the 5'-Cap and insertion of a small synthetic oligo. This sequence could also show up in the sequencing and be shown in the genbank, but removed when aligned to the genome in the UCSC browser.

*Source: 1. Suzuki, Y., Yoshitomo-Nakagawa, K., Maruyama, K., Suyama, A. & Sugano, S. Construction and characterization of a full length-enriched and a 5'-end-enriched cDNA library. Gene 200, 149-156 (1997).*

### Question 2: ERA and ERB

Our collaborators designed a ChIP study using so-called tiling arrays (an outdated technique these days, but the top of the pop at the time: see [http://en.wikipedia.org/wiki/Tiling\\_array](http://en.wikipedia.org/wiki/Tiling_array)): one for estrogen receptor alpha (ERA), one for estrogen receptor beta (ERB). All the sites are stored in BED files respectively for two ERs. These are now available in the homework directory, and are both mapped on hg18 genome. The

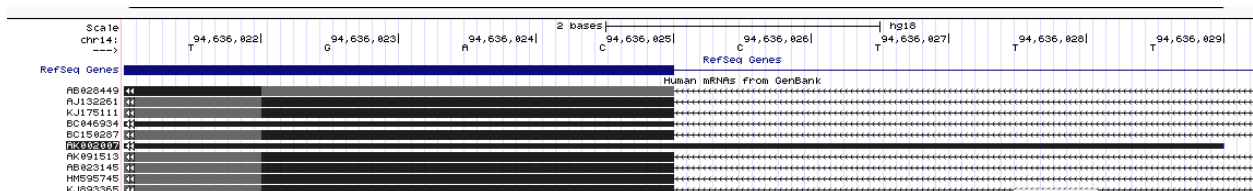


Figure 1: UCSC screenshot, showing the five first bases of transcription of AK002007

```

    endonuclease Dicer (EC 3.1.26.-)
    /codon_start=1
    /protein_id="BAG51002.1"
    /translation="MVVSIFDPPVNWLPFGYVVNQDKSNTDKWEKDEMTKDCMLANGK
    LDEDYEEDEEEESLMWRAPKEEADYEDDFLEYDQEHIRFIDNMLMGSGAFVKKISLS
    PFSTTDSAYEWKMPKKSSSLGSMFSSDFEDFDYSSWDAMCYLDPSKAVEEDDFVVGFW
    NPSEENCVDGTGQSI SYDLHTEQCIADKSIADCEALLGCYLTSCGERAAQLFLCSL
    GLKVLFPVIKRTDREKALCPTRENFNSQQKNLSVSCAAASVASSRSSVLKDSEYGLKI
    PPRCMFDHPDADKTLNHLISGFENFEKKINRYRFKNKAYLLQAFTHASYHYNTITDCYQ
    RLEFLGDAILDYLITKHLIEDPRQHSFPGVLTDLRSALVNNTIFASLAVKYDYHKYFKA
    VSPELFHVIDDVFQFQLEKNEMQGMDELRRSEEDKEEDIEVPKAMGDIFESLAGA
    IYMDSGMSLETWQVYYPMMRPLIEKFSANVPRSPVRELLEMEPETAKFSPAERTYDG
    KVRVTVEVVGKGFKGVGRSYRIAKSAAARRALRSLKANQPQVNS"

ORIGIN
1  gaagcaaaaa ggtcagcaac tgtaatctgt atcgcccttg aaaaaagaag ggactaccca
61  gccgcatggt ggtgtcaata tttgatcccc ctgtgaattg gcttcctcct gggtatgtag
121  taaatcaaga caaaagcaac acagataaat gggaaaaaga tgaatgaca aaagactgca
181  tgctggcgaa tggcaaaact gatgaggatt acgaggagga ggatgaggag gaggagagcc
241  tgatgtggag ggctccgaag gaagaggctg actatgaaga tgatttcctg gagtatgac
301  aggaacacat cagatttata gataatatgt taatggggtc aggagctttt gtaaaagaaa
361  tctctcttcc tctcttttca accactgatt ctgcatatga atggaaaatg cccaaaaaat
421  cctccttagg tagtatgcca ttttcatcag attttgagga ttttgactac agctcttggg
481  atgcaatgtg ctatctggat cctagcaaaag ctgttgaaga agatgacttt gtggtggggt
541  tctggaatcc atcagaagaa aactgtggtg ttgacacggg aaagcagtc cttctttacg
601  acttgcacac tgagcagtggt attgtcgaca aaagcatagc ggactgtgtg gaagccctgc
661  tgggctgcta tttaaccagc tgtggggaga gggctgctca gcttttcctc tgttcactgg
721  ggctgaaggt gctcccggtt attaaaagga ctgacgggga aaaggccctg tgccctactc
781  gggagaattt caacagccaa caaaagaacc tttcagtgag ctgtgctgct gcttctgtgg
841  ccagttcacg ctctctgtga ttgaaagact cggaatatgg ttgtttgaag attccaccaa
901  gatgtatgtt tgatcatcca gatgcagata aaacactgaa tcaccttata tcgggggttg
961  aaaattttga aaagaaaatc aactacagat tcaagaataa ggcttacctt ctccaggctt
1021  ttacacatgc ctctaccac tacaatacta tcaactgattg ttaccagcgc ttagaattcc
1081  tgggagatgc gattttggac tacctcataa ccaagcacct ttatgaagac ccgcggcagc
1141  actccccggg ggtcctgaca gacctgcggt ctgccctggt caacaacacc atctttgcat
1201  cgctggctgt aaagtacgac taccacaagt acttcaaagc tgtctctcct gagctcttcc
1261  atgtcattga tgactttgtg cagtttcagc ttgagaagaa tgaaatgcaa ggaatggatt
1321  ctgagcttag gagatctgag gaggatgaag agaaagaaga ggatattgaa gttccaaagg
1381  ccattggggg tatttttgag tcgcttgctg gtgccattta catggatagt gggatgtcac
1441  tggagacagt ctggcaggtg tactatccca tgatgcggcc actaatagaa aagttttctg
1501  caaatgtacc ccgttcccct gtgcgagaat tgcttgaagt ggaaccagaa actgccaaat
1561  ttagcccggc tgagagaact tacgacggga aggtcagagt cactgtggaa gtagtaggaa
1621  aggggaaatt taaaggtgtt ggtcgaagtt acaggattgc caaatctgca gcagcaagaa
1681  gagccctccg aagcctcaaa gctaataaac ctacaggctcc caatagctga aaccgctttt
1741  taaaattcaa aacaagaaac

```

Figure 2: Screenshot from the GenBank database of AK002007. The 7 nucleotides which are different between the UCSC browser and the GenBank entry are highlighted in blue.

current situation is that we know to some degree what ERA does, but not what ERB does (there are some evidence that they share some functions, but not all). So, we need bigger experiments and better statistics.

- a) Using BEDtools within Linux: What is the genome coverage (% of base pair covered at each chromosome) for ERB and ERA sites? If you need a file with chromosome sizes for hg18, it included in the assignment: hg18\_chrom\_sizes.txt. Plot the fractions for all chromosomes as a single barplot in R. Briefly comment the results. Is there anything particularly surprising? Try to explain the outcome (biological and/or experimental setup explanations)?
- b) Again, using BEDtools in Linux: How many ERA sites do/do not overlap ERB sites, and vice versa? Show the Linux commands and then a Venn diagram summarizing the results. The Venn diagram can be made in R using one of many venn diagram packages, but you can also make it in any drawing program.

### **Question 3: Ribosomal Gene (\*)**

Your group just got this email from a frustrated fellow student:

My supervisor has found something he thinks is a new ribosomal protein gene in mouse. It is at chr9:24,851,809-24,851,889, assembly mm8. His arguments for this are a) It has high conservation in other species because ribosomal protein genes from other species map to this mouse region b) And they are all called Rpl41 in the other species (if you turn on the other Refseq you see this clearly in fly and other species).

But, I found out that if you take the fly refseq sequence mentioned above (from Genbank) and BLAT this to the fly genome, you actually get something that looks quite different from the one in the mouse genome. How can this be? Is the mouse gene likely to be real? If not, why? (Maximum 20 lines, plus possibly genome browser pictures)

**Answer:**