

# Alternative CDF environments

Laurent Gautier

November 15, 2007

## Introduction

On short oligonucleotide arrays, several probes are designed to match a target transcript, and probes matching the same target transcript can be grouped in a probe set. Between the time the probes for a given short oligonucleotide chip were designed, and the time an analysis is made, the knowledge of expected transcripts for a given organism might have changed. Unless one includes the latest development in transcripts into an analysis, the analysis could suffer from what we like to call a *Dorian Gray*<sup>1</sup> effect. The chip itself does not change, which means that the probes and their respective sequences remain the same, while the knowledge of the transcripts, and eventually their sequence, might evolve, and in time the immobility of the probe and probe sets give an uglier picture of the biological phenomena to study. Being able to easily modify or replace the grouping of probes in probe sets gives the opportunity to minimize this effect.

The package is directly usable with *Affymetrix GeneChip* short oligonucleotide arrays, and can be adapted or extended to other platforms.

The bibliographic reference associated with the package is:

Alternative mapping of probes to genes for *Affymetrix* chips Laurent Gautier, Morten Mooller, Lennart Friis-Hansen, Steen Knudsen BMC Bioinformatics 2004, 5:111 (14–August 2004)<sup>2</sup>.

Let's start by loading the package:

```
> library(altcdfenvs)
```

## The class CdfEnvAffy

Each instance of this class contains a way to group probes in probe sets. Different instances, describing different ways to group probes in probe sets, can co-exist for a given chip type.

---

<sup>1</sup>From the novel 'The Picture of Dorian Gray' by Oscar Wilde.

<sup>2</sup>An other research team mapped probes to NCBI's RefSeq at about the same time that we were doing it independantly. For details, see: Mecham *et al.*, 2004, Nucl. Acids Research, vol. 32, no. 9, pp. e74-e74.

When experimenting, it is highly recommended to use the functions `validCdfEnvAffy` and `validAffyBatch` to make sure that the instance fulfills basic requirements for consistency.

## Reading sequence information in FASTA connections

The package contains simple functions to read **R** connections in the FASTA format. Typically, collections of sequences are stored in FASTA files, which can be significantly large, one can wish to read and process sequences one after the other. This can be done by opening the file in 'r' mode:

```
> fasta.filename <- system.file("exampleData", "sample.fasta",  
+   package = "altcdfenvs")  
> con <- file(fasta.filename, open = "r")
```

Reading the sequences one after another, and printing information about them in turn goes like:

```
> fasta.seq <- read.FASTA.entry(con)  
> while (!is.null(fasta.seq$header)) {  
+   print(fasta.seq)  
+   fasta.seq <- read.FASTA.entry(con)  
+ }
```

FASTA sequence:

```
>gnl|UG|Hs#S1730546 membrane-spanning 4-domains, subfamily A ...  
AACCCATTTCAACTGCCTATTCAGAGCATGCAGTAAGAGGAAATCCACCAAGTCTCAATA ...
```

FASTA sequence:

```
>gi|28626515|ref|NM_007257.3| Homo sapiens paraneoplastic an ...  
GGTCATTTGTCCAGAAAACCTTTGTGACTGTCTTTGAGTGACCTAGTCTGGGACCCATTCA ...
```

FASTA sequence:

```
>gi|31377729|ref|NM_020143.2| Homo sapiens putatative 28 kDa ...  
TGGCTTCTGCGTGGTGCAGCTGCGCACGTGTTTCAGCCGGCAGCGCTTTAAGATTCCGG ...
```

```
> close(con)
```

One can foresee that the matching of a set of reference sequences against all the probes can be parallelized easily: the reference sequences can simply be distributed across different processors/machines. When working with all the reference sequences in a single large FASTA file, the option `skip` can let one implement a poor man's parallel sequence matching very easily.

# Creating an alternative mapping from sequences in a FASTA file

## Select the constituting elements

- Chip type: For this tutorial we decide to work with the Affymetrix chip HG-U133A.
- Target sequences: The set of target sequences we use for this tutorial is in the exemplar FASTA file:

```
> con <- file(fasta.filename, open = "r")
> n <- countskip.FASTA.entries(con)
> close(con)
> con <- file(fasta.filename, open = "r")
> my.entries <- read.n.FASTA.entries.split(con, n)
> close(con)
```

## matching the probes

The package *matchprobes* and the probe data package for HG-U133A are required to perform the matching. The first step is to load them:

```
> library(matchprobes)
> library(hgu133aprobe)
```

The matching is done simply (one can refer to the documentation for the package *matchprobes* for further details):

```
> m <- matchprobes(my.entries$sequences, hgu133aprobe$sequence,
+   probepos = TRUE)
```

The package contains a function to create a *CdfEnv* from the matches:

```
> get.RNA.IDs <- function(x) {
+   reg <- regexpr("(Hs#|NM)[^[:blank:]]+", x)
+   r <- substr(my.entries$headers, reg, reg + attr(reg, "match.length") -
+     1)
+   return(r)
+ }
> ids <- get.RNA.IDs(my.entries$headers)
> alt.cdf <- buildCdfEnv.matchprobes(m, ids, nrow.chip = 712, ncol.chip = 712,
+   chiptype = "HG-U133A", probes.pack = "hgu133aprobe")
```

Note that the size for chip must be specified. This is currently a problem with *cdfenvs* as they are created by the package *makecdfenv*. The class *CdfEnv* suggests a way to

solve this (hopefully this will be integrated in *makecdfenv* in the near future). When this happens, the section below will be replaced by something more intuitive. But in the meanwhile, here is the current way to use our shiny brand new `CdfEnv`:

```
## say we have an AffyBatch of HG-U133A chips called 'abatch'

## summary checks to avoid silly mistakes
validAffyBatch(abatch, alt.cdf)

## it is ok, so we proceed...

## get the environment out of it class
alt.cdfenv <- alt.cdf@envir

abatch@cdfName <- "alt.cdfenv"
```

From now on, the object `abatch` will use our ‘alternative mapping’ rather than the one provided by the manufacturer of the chip:

```
print(abatch)
```