

The Art of Multiple Sequence Alignment in R

Erik S. Wright
University of Wisconsin
Madison, WI

October 7, 2014

Contents

1	Introduction	1
2	Alignment Speed	2
3	Alignment Accuracy	3
4	Example: Single Gene Alignment	4
5	Example: Building a Guide Tree	4
6	Session Information	6

1 Introduction

This document is intended to illustrate the art of multiple sequence alignment in *R* using DECIPHER. Even though its beauty is often concealed, multiple sequence alignment is a form of art in more ways than one. Take a look at Figure 1 for an illustration of what is happening behind the scenes during multiple sequence alignment. The practice of sequence alignment is one that requires a degree of skill, and it is that art which this vignette intends to convey. It is not simply enough to “plug” sequences into a multiple sequence aligner and blindly trust the result. An appreciation for the art as well a careful consideration of the results are required.

What really is multiple sequence alignment, and is there a single correct alignment? Generally speaking, alignment seeks to perform the act of taking multiple divergent biological sequences of the same “type” and fitting them to a form that reflects some shared quality. That quality may be how they look structurally, how they evolved from a common ancestor, or optimization of a mathematical construct. As with most multiple sequence aligners, DECIPHER

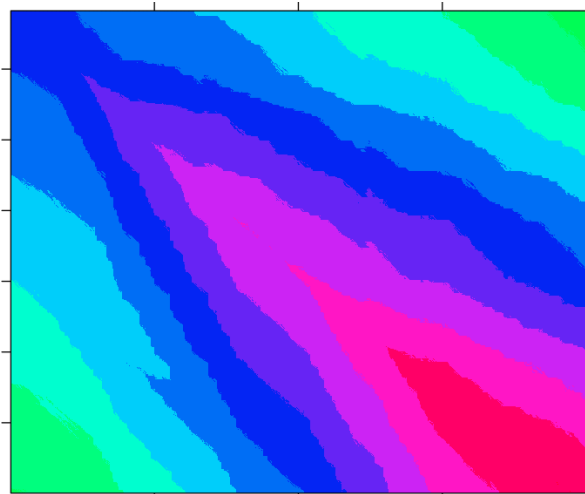


Figure 1: The art of multiple sequence alignment.

is “trained” to maximize scoring metrics in order to accomplish a combination of both structural alignment and evolutionary alignment. The idea is to give the alignment a biological basis even though the sequences will never meet each other and align under any natural circumstance.

The workhorse for sequence alignment in DECIPHER is `AlignProfiles`, which takes in two aligned sets of DNA, RNA, or amino acid sequences and returns a merged alignment. For more than two sequences, the function `AlignSeqs` will perform multiple sequence alignment in a progressive/iterative manner on sequences of the same kind. In this case, multiple alignment works by aligning two sequences, merging with another sequences, merging with another alignment, and so-forth until all the sequences are aligned. This process is iterated to further refine the alignment. There are other functions that extend use of `AlignSeqs` for different purposes:

1. The first is `AlignTranslation`, which will align DNA/RNA sequences based on their amino acid translation and then reverse translate them back to DNA/RNA. This method may improve both alignment accuracy and speed, since amino acid sequences are more conserved, have a well-studied structural basis, and are shorter than their corresponding coding sequences.
2. The second function, `AlignDB`, enables generating alignments from many more sequences than possible to fit in memory. Its main purpose is to merge sub-alignments where each alignment alone is composed of many sequences. This is accomplished by storing all of the sequences in a database and only working with “profiles” representing the sequences.

2 Alignment Speed

The dynamic programming method for aligning two profiles requires order $N \times M$ time and memory space where N and M are the width of the pattern and subject. Since multiple sequence alignment is an inherently challenging problem for large sequences, heuristics are employed to maximize speed while maintaining reasonable accuracy. In this regard, the two most important parameters accessible by the user are *restrict* and *anchor*. The objective of the *restrict* parameter is to convert the problem from one taking quadratic time to linear time. The goal of the *anchor* parameter is do the equivalent for memory space so that very long sequences can be efficiently aligned.

The orange anti-diagonal line in Figure 2 shows the optimal path for aligning two sequence profiles. The blue segments to the left and right of the optimal path give the constraint boundaries, which the user controls with the *restrict* parameter. Areas above and below the upper and lower constraint boundaries are neglected from further consideration. A higher (less negative) value of *restrict* will further constrain the possible “alignment space,” which represents all possible alignments between two sequences. Since the optimal path is not known till completion of the matrix, it is risky to overly constrain the matrix. This is particularly true in situations where the sequences are not mostly overlapping because the optimal path will likely not be diagonal, causing the path to cross a constraint boundary. In the non-overlapping case *restrict* should be set below the default to ensure that the entire “alignment space” is available.

Neglecting the “corners” of the alignment space effectively converts a quadratic time problem into a near-linear time problem. We can see this by comparing `AlignProfiles` with and without restricting the matrix at different sequence lengths. To extend our comparison we can include the `Biostrings` function

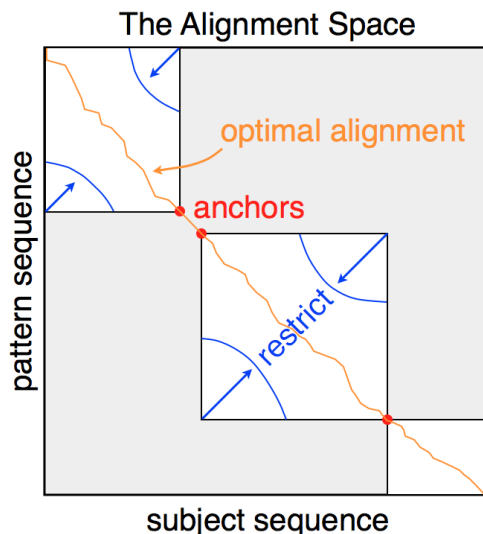


Figure 2: The possible alignment space.

`pairwiseAlignment`. In this simulation, two sequences with 90% identity are aligned and the elapsed time is recorded for a variety of sequence lengths. As can be seen in Figure 4 below, *without* restriction `AlignProfiles` takes quadratic time in the same manner as `pairwiseAlignment`. However, *with* restriction `AlignProfiles` takes linear time, taking only a few microseconds per nucleotide.

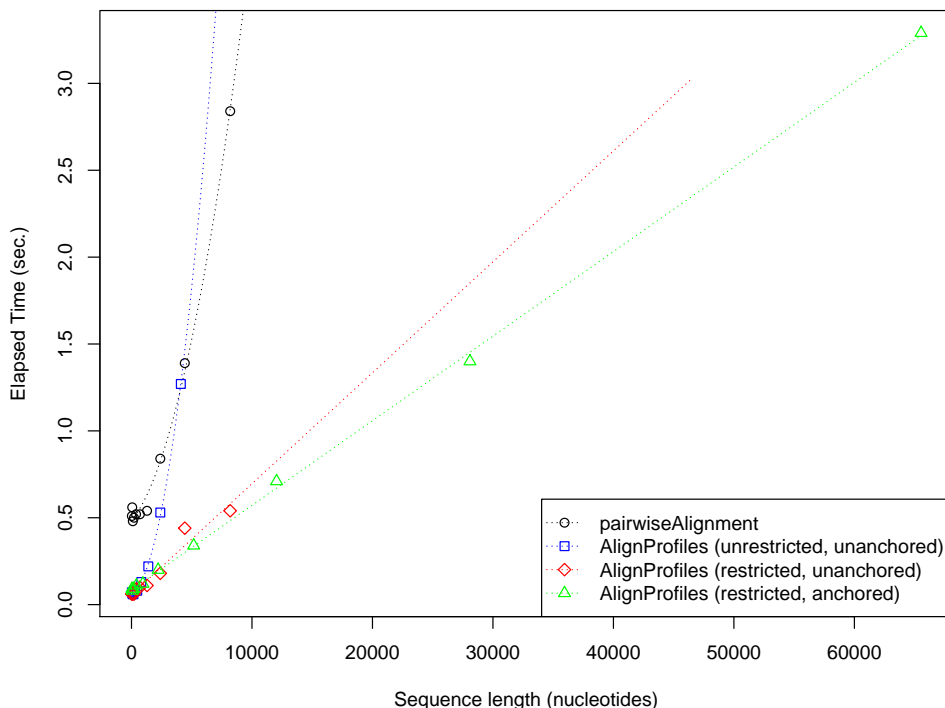


Figure 3: Global Pairwise Sequence Alignment Timings.

The parameter *anchor* controls the fraction of sequences that must share a common region to anchor the alignment space (Fig. 2). `AlignProfiles` will search for shared anchor points between the two sequence sets being aligned, and if the fraction shared is above *anchor* then that position is fixed in the “alignment space.” Anchors are 15-mer (for DNA/RNA) or 7-mer (for AA) exact matches between two sequences that must occur in the same order as they do in the sequences. Anchoring generally does not affect accuracy, but can greatly diminish the amount of memory required for alignment. The longest pair of sequence profiles that can be aligned without anchoring is 46 thousand nucleotides as shown by the end of the red dashed line in Figure 3. If regularly spaced anchor points are available then the maximum sequence length is greatly extended. However, anchoring comes with the caveat that it can in some circumstances result in a different alignment, sometimes better and sometimes worse than without anchoring.

3 Alignment Accuracy

Figure 4 compares the performance of DECIPHER and other sequence alignment software on structural amino acid benchmarks [1]. All structural benchmarks have flaws, some of which can easily be found by eye in highly similar sequence sets, and therefore benchmark results should be treated with care [3]. As can be seen in the figure, the performance of DECIPHER is similar to that of other popular alignment software such as

ClustalW [6]. Most importantly, the accuracy of amino acid alignment begins to drop-off when sequences in the reference alignment have less than 40% average pairwise identity. A similar decline in performance is observed with DNA/RNA sequences, but the drop-off occurs much earlier at around 60% sequence identity. Therefore, it is generally preferable to align coding sequences by their translation using `AlignTranslation`. This function first translates the input DNA/RNA sequences, then aligns the translation, and finally reverse translates to obtain aligned DNA/RNA sequences.

4 Example: Single Gene Alignment

For this example we are going to align the *rplB* coding sequence from many different Bacteria. The *rplB* gene encodes one of the primary rRNA binding proteins: the 50S ribosomal protein L2. We begin by loading the library and importing the sequences from a FASTA file:

```
> library(DECIPHER)
> # specify the path to your sequence file:
> fas <- "<<path to FASTA file>>"
> # OR find the example sequence file used in this tutorial:
> fas <- system.file("extdata", "50S_ribosomal_protein_L2.fas", package="DECIPHER")
> dna <- readDNAStringSet(fas)
> dna # the unaligned sequences
A DNAStringSet instance of length 317
      width seq                                     names
[1]   819 ATGGCTTTAAAA...AAAAAAGAAAA Rickettsia prowaz...
[2]   822 ATGGGAATACGT...AGAAGGAAAAAG Porphyromonas gin...
[3]   822 ATGGGAATACGT...AGAAGGAAAAAG Porphyromonas gin...
[4]   822 ATGGGAATACGT...AGAAGGAAAAAG Porphyromonas gin...
[5]   819 ATGGCTATCGTT...CGTCGTGGTAAA Pasteurella multo...
...   ...
[313] 819 ATGGCAATTGTT...CGCCGTACTAAA Pectobacterium at...
[314] 822 ATGCCTATTCAA...CGTCGCGTCAAG Acinetobacter sp....
[315] 864 ATGGGCATTTCGC...GGTCGTCACTCT Thermosynechococc...
[316] 831 ATGGCACTGAAG...AAGCGGAAGAAG Bradyrhizobium ja...
[317] 840 ATGGGCATTTCGC...TCCGGGCGAGGT Gloeobacter viola...
```

We can align the DNA by either aligning the coding sequences or their translations (amino acid sequences). Both methods result in an aligned set of DNA sequences, unless the argument *asAAStringSet* is `TRUE` in `AlignTranslation`. A quick inspection reveals that the method of translating before alignment yields a more appealing result. However, if the dna did not belong to a coding sequence then the only option would be to use `AlignSeqs`.

```
> AA <- AlignTranslation(dna, asAAStringSet=TRUE) # align the translation
> BrowseSequences(AA, highlight=1) # view the alignment
> DNA <- AlignTranslation(dna) # align the translation then reverse translate
> DNA <- AlignSeqs(dna) # align the sequences directly without translation
```

5 Example: Building a Guide Tree

The `AlignSeqs` uses a guide tree to decide which order to align pairs of sequence profiles. The *guideTree* input is a *data.frame* with the grouping of each sequence at increasing levels of dissimilarity between the groups. By default this guide tree is generated directly from the sequences using the order of shared k-mers

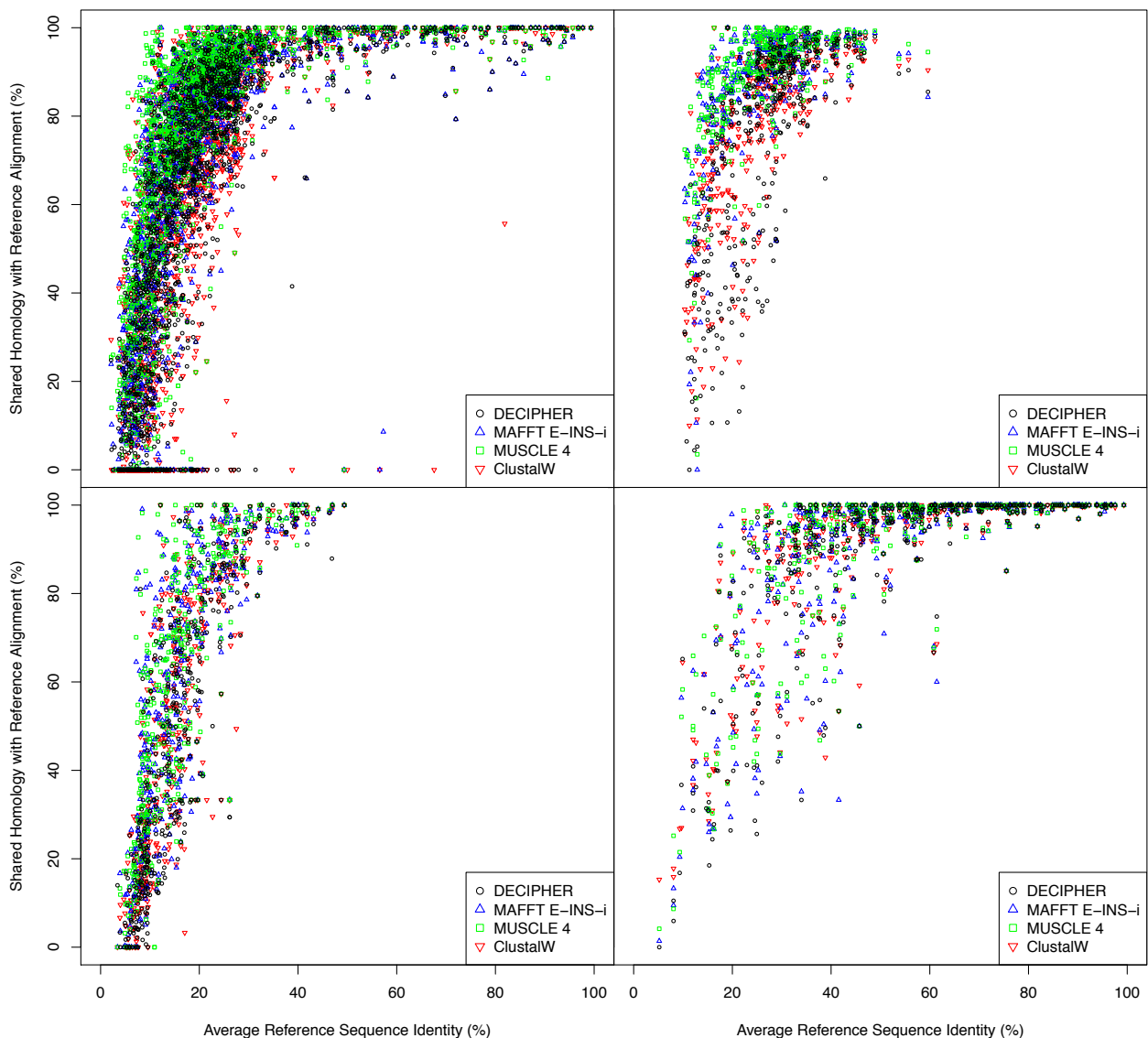


Figure 4: Performance comparison between different programs for multiple alignment using amino acid structural benchmarks ([6], [2], [4]). The x-axis shows percent identity between sequences in each reference alignment. The y-axis gives the percentage of correctly aligned residues in the estimated alignment according to the reference alignment (SP-score). The upper-left plot is for the PREFAB (version 4) benchmark [2]. The upper-right plot shows the results of the BALIBASE (version 3) benchmark [7]. The lower-left plot is for SABMARK (version 1.65) [8]. The lower-right plot gives the results on the OXBENCH alignments [5].

(the argument *guideTree* is `NULL`). However, in some circumstances it may be desirable to provide a guide tree as input.

The first circumstance in which a guide tree could be useful is when the sequences are highly divergent. In such a case the quick guide tree generated by `AlignSeqs` may perform worse than a more accurate guide tree generated from pairwise alignments.

```
> # form guide tree using pairwiseAlignment
> l <- length(dna)
> d <- matrix(0, nrow=l, ncol=l)
> pBar <- txtProgressBar(style=3)
> for (j in 2:l) {
  d[j, 1:(j - 1)] <- pairwiseAlignment(rep(dna[j], j - 1),
    dna[1:(j - 1)],
    scoreOnly=TRUE)
  setTxtProgressBar(pBar, j/l)
}
> close(pBar)
> # rescale the distance scores from 0 to 1
> m <- max(d)
> d[lower.tri(d)] <- d[lower.tri(d)] - m
> m <- min(d)
> d[lower.tri(d)] <- d[lower.tri(d)]/m
> # form a guide tree from the distance matrix
> gT <- IdClusters(d, cutoff=seq(0, 1, 0.01))
> # use the guide tree as input for alignment
> DNA <- AlignSeqs(dna, guideTree=gT) # align directly
> DNA <- AlignTranslation(dna, guideTree=gT) # align by translation
```

A second circumstance is if there are a large number of sequences (more than 46,340), in which case an alternative guide tree method is required. In this case a rough guide tree can be generated by directly clustering the sequences.

```
> # form guide tree using inexact clustering
> gT <- IdClusters(myXStringSet=dna, method="inexact", cutoff=seq(0.05, 0.9, 0.05))
> # use the guide tree as input for alignment
> DNA <- AlignSeqs(dna, guideTree=gT) # align directly
> DNA <- AlignTranslation(dna, guideTree=gT) # align by translation
```

6 Session Information

All of the output in this vignette was produced under the following conditions:

- R version 3.1.1 (2014-07-10), i386-w64-mingw32
- Base packages: base, datasets, grDevices, graphics, methods, parallel, stats, utils
- Other packages: BiocGenerics 0.10.0, Biostrings 2.32.1, DBI 0.3.1, DECIPHER 1.10.1, IRanges 1.22.10, RSQLite 0.11.4, XVector 0.4.0
- Loaded via a namespace (and not attached): stats4 3.1.1, tools 3.1.1, zlibbioc 1.10.0

References

- [1] Edgar, R. C. Quality measures for protein alignment benchmarks. *Nucleic Acids Research*, 38(7), 2145-2153. doi:10.1093/nar/gkp1196, 2010.
- [2] Edgar, R. C. MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids Research*, 32(5), 1792-97, 2004.
- [3] Iantorno, S., Gori, K., Goldman, N., Gil, M., & Dessimoz, C. Who watches the watchmen? An appraisal of benchmarks for multiple sequence alignment. *Methods in Molecular Biology (Clifton, N.J.)*, 1079, 59-73. doi:10.1007/978-1-62703-646-7_4, 2014.
- [4] Katoh, K., Misawa, K., Kuma, K.-I., & Miyata, T. MAFFT: a novel method for rapid multiple sequence alignment based on fast Fourier transform. *Nucleic Acids Research*, 30(14), 3059-3066, 2002.
- [5] Raghava, G. P., Searle, S. M., Audley, P. C., Barber, J. D., & Barton, G. J. OXBench: a benchmark for evaluation of protein multiple sequence alignment accuracy. *BMC Bioinformatics*, 4: 47, 2003.
- [6] Thompson, J. D., Higgins, D. G., & Gibson, T. J. CLUSTAL W: improving the sensitivity of progressive multiple sequence alignment through sequence weighting, position-specific gap penalties and weight matrix choice. *Nucleic Acids Research*, 22(22), 4673-4680, 1994.
- [7] Thompson, J. D., Koehl, P., Ripp, R., & Poch, O. BALiBASE 3.0: Latest developments of the multiple sequence alignment benchmark. *Proteins*, 61(1), 127-136, 2005.
- [8] Van Walle, I., Lasters, I., & Wyns, L. SABmark—a benchmark for sequence alignment that covers the entire known fold space. *Bioinformatics*, 21(7), 1267-1268, 2005.