



HAL
open science

A Survey on Algorithmic Aspects of Tandem Repeats Evolution

Eric Rivals

► **To cite this version:**

Eric Rivals. A Survey on Algorithmic Aspects of Tandem Repeats Evolution. International Journal of Foundations of Computer Science, 2004, 15 (2), pp.225-257. 10.1142/S012905410400239X . lirmm-00108543

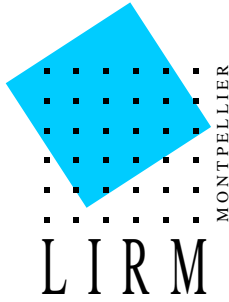
HAL Id: lirmm-00108543

<https://hal-lirmm.ccsd.cnrs.fr/lirmm-00108543>

Submitted on 23 Oct 2006

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



LABORATOIRE D'INFORMATIQUE, DE
ROBOTIQUE ET DE MICROÉLECTRONIQUE DE
MONTPELLIER

Unité Mixte CNRS - Université Montpellier II C 5506

RAPPORT DE RECHERCHE

A Survey on Algorithmic Aspects of Tandem Repeats Evolution

Eric Rivals

20/08/2003

R.R.LIRMM 03-017

A Survey on Algorithmic Aspects of Tandem Repeats Evolution

Auteurs : Eric Rivals

Mèls / Emails : rivals@lirmm.fr

Résumé : Les répétitions locales dans les génomes sont appelées *répétitions en tandem*. Elles sont constituées de plusieurs copies légèrement différentes d'un motif répété et changent durant l'évolution lorsque ces copies subissent des mutations ponctuelles, lorsque des événements d'*amplification* ajoutent de nouvelles copies par duplication ou des copies sont enlevées par *contraction* de copies identiques. Ces modifications font évoluer les répétitions en tandem de manière dynamique. Cette constatation soulève deux problèmes : la *Reconstruction d'une Histoire de Duplications* vise à retrouver l'histoire des amplifications et des mutations qui ont produit une répétition donnée à partir d'une copie ancestrale du motif. Étant données deux répétitions en tandem prélevées à une même position génomique dans deux individus différents, nommées *allèles*, et une fonction de coût pour les amplifications, contractions et mutations, l'*Alignement d'Allèles* consiste à trouver l'alignement des deux séquences de coût minimal. Nous présentons ici un état de l'art sur ces deux problèmes qui permettent d'étudier les mécanismes évolutifs propres aux répétitions en tandem. La présentation est étayée d'un survol des applications biologiques dans lesquelles les répétitions en tandem interviennent.

Abstract: Local repetitions in genomes are called *tandem repeats*. A tandem repeat contains multiple, but slightly different copies of a repeated unit. It changes over time as the copies are altered by mutations, when additional copies are created by *amplification* of an existing copy, or when a copy is removed by *contraction*. These changes let tandem repeats evolve dynamically. From this statement follow two problems. TANDEM REPEAT HISTORY aims at recovering the history of amplifications and mutations that produced the tandem repeat sequence given as input. Given the tandem repeat sequences at the same genomic location in two individuals and a cost function for amplifications, contractions, and mutations, the purpose of TANDEM REPEAT ALLELE ALIGNMENT is to find an alignment of the sequences having minimal cost. We present a survey of these two problems that allow to investigate evolutionary mechanisms at work in tandem repeats.

Mots clés : mots, périodicité, répétition en tandem, évolution, duplication, phylogénie, revue, alignement.

Keywords: periodic words, tandem repeat, evolution, duplication, phylogeny, survey, alignment.

A Survey on Algorithmic Aspects of Tandem Repeats Evolution

Eric Rivals

L.I.R.M.M., CNRS U.M.R. 5506
161 rue Ada, F-34392 Montpellier Cedex 5, France
rivals@lirmm.fr

Abstract. Local repetitions in genomes are called *tandem repeats*. A tandem repeat contains multiple, but slightly different copies of a repeated unit. It changes over time as the copies are altered by mutations, when additional copies are created by *amplification* of an existing copy, or when a copy is removed by *contraction*. These changes let tandem repeats evolve dynamically. From this statement follow two problems. TANDEM REPEAT HISTORY aims at recovering the history of amplifications and mutations that produced the tandem repeat sequence given as input. Given the tandem repeat sequences at the same genomic location in two individuals and a cost function for amplifications, contractions, and mutations, the purpose of TANDEM REPEAT ALLELE ALIGNMENT is to find an alignment of the sequences having minimal cost. We present a survey of these two problems that allow to investigate evolutionary mechanisms at work in tandem repeats.

1 Introduction

A striking genetic difference between species is the size of their genome. Relatively simple organisms, like the protist *Amoeba dubia*, may have much larger genome than *Homo sapiens* for instance. These dramatic differences are due to the presence of repeats. In general, in eukaryotes, organisms whose cells bear a kernel, duplicated genetic material is abundant and can account for up to 60% of the genome. Although some of the mechanisms that generate these repeats are known, from the point of view of evolution, the reasons for such redundancy remain an enigma.

Repeats whose copies are distant in the genome, whether or not located on the same chromosome, are called distant repeats. In this review, we focus on repeats whose copies are adjacent on a chromosome. Because of this characteristic, they bear the name of **tandem repeats**. Among those, biologists distinguish **micro-satellites**, **mini-satellites**, and **satellites**, according to the length of their repeated unit: between 1 and 6 base-pairs, between 7 and 50 base-pairs¹, and above 50 base-pairs¹, respectively. These names are mainly used for repeats located in regions that do not contain genes. In addition to these sub-classes, numerous groups of similar genes that originate from the same ancestor gene are organized in tandem. They are termed **tandemly repeated genes**.

Local repeats in the DNA arise, grow or disappear through molecular events that copy a contiguous segment on the DNA and insert one or many copies of it next to the original segment, or perform the dual operation. We name these two types of events **amplification** and **contraction**. Like any other segment of the genome, the repeated copies also change through **point mutations**: insertion, deletion or substitution of one base. Point mutations give rise to approximate tandem repeats. The pattern of point mutations along the tandem array of copies informs us on the parent-child relationships between copies. In other words, it gives access to the history of the tandem repeat.

The relatively high frequency of these events let these local repeats evolve rapidly. For a given species and at a precise location on the chromosome, a **locus**, the repeat varies in sequence and/or length in different individuals. Hence, such a locus is said to be **polymorphic** and each different sequence encountered at this locus is called an **allele**.

¹ Chromosomes are made of a double-stranded Deoxyribonucleic Acid (DNA) helix, whose basic building block is a pair of bases. The unit of a DNA sequence is thus called a **base-pair** and is abbreviated by bp.

1.1 Approximate Tandem Repeats

In biology, local repetitions in DNA are called "tandem repeats" irrespectively of the number of copies. In computer science, a local repetition is dubbed a **square** if it contains two copies, a **cube** if it contains three, and so on.

An amplification creates a substring that is an **Exact Tandem Repeat**, ETR for short. An ETR is a power of the original pattern: for an integer m , it equals u^m if the pattern is u . When later in the course of evolution point mutations affect this ETR, they let identical positions in adjacent copies differ and the ETR becomes an **Approximate Tandem Repeat**, ATR for short. Note that any sequence is an ATR of some motif. In practice, only repeats whose copies are similar enough receive attention. The level of internal similarity that distinguishes any random sequence from a sequence of true repeats, i.e., that is created by some amplifications, is defined from a statistical view-point (for example in the software TRF [Ben99]) or by an information theoretical measure ([RDDD96,RDD+97]). The problem of detecting significant ETR or ATR is an active area of research (see for instance [RDD+97,SM98,DDR99,Ben99,KK00,KK01,SG02]). In the sequel of the paper, by ATR we mean a tandem repeat with sufficient internal similarity. An example of an ATR is given Fig. 1 under the form a multiple alignment of its copies.

Point mutations could cause two adjacent copies to diverge so far that their common ancestry is not recognizable anymore from sequence similarity. In this case, it is not a repeat anymore. A major hypothesis is that amplification is favored by the similarity of adjacent patterns, and that when copies have diverged for a long time such former repeat does not undergo amplification anymore. In highly polymorphic loci, like some minisatellites, amplifications and contractions are more probable than point mutations. On the contrary, tandemly repeated genes can accumulate hundreds of mutations and still undergo some amplifications; in this case, amplifications and contractions are less frequent than point mutations.

When one wishes to establish the common ancestry of any two genes, one first searches for sequence similarity. The similarity is quantified through sequence alignment. The ALIGNMENT is a weighted version of the LONGEST COMMON SUBSEQUENCE problem and, in the classical setup, considers only point mutations. An exact solution is based on dynamic programming [Gus97,SK99]. Dealing with tandem repeat requires to consider also amplifications and contractions. We do not report on other algorithmic and combinatorial problems on local repetitions and refer the reader to numerous textbooks on the subject, among which [Lot99,CHL01,Gus97].

c	t	g	a	g	c	t	c	A	a	C	c	t	t	g	c	t	c	T	g	a	g	c	A	T	c	a	t	c	t	t	-	c	t
c	t	g	a	g	c	t	c	c	a	t	c	t	t	A	c	A	c	T	g	a	g	A	A	G	c	a	C	c	t	G	-	c	t
G	C	A	a	g	c	t	c	c	a	t	c	t	t	g	c	t	T	G	g	a	g	c	t	c	c	T	t	c	t	t	g	c	t
c	C	A	a	g	c	t	c	T	a	t	c	-	t	A	c	t	c	c	A	a	g	c	t	c	c	a	t	c	t	t	g	c	t
c	A	g	a	g	c	t	c	c	a	t	c	-	t	g	c	t	c	c	A	a	g	c	t	c	c	a	t	c	t	t	g	c	t
c	G	A	a	g	T	G	c	c	a	-	A	t	C	g	c	t	c	c	A	a	g	c	A	c	T	a	t	c	t	t	g	c	t
G	t	g	a	g	c	A	A	c	a	t	c	-	t	g	c	A	T	A	g	a	C	A	t	T	c	a	t	c	t	t	a	c	t
c	A	g	a	g	c	t	c	c	a	t	c	t	A	g	-	t	c	A	g	a	g	A	t	c	c	a	t	c	C	A	-	c	t

Fig. 1. A multiple alignment of the 8 copies of a tandem repeat found on the human chromosome 22. The lines of the alignment contain the copy in the same order than on the chromosome. Symbols in bold uppercase mark differences between the current copy and a 34 bp consensus motif. On the third column from the right, the copies 3 to 6 all have an extra *g* character suggesting that they may have arisen through an amplification of arity 4 after the *g* was inserted in the original copy.

1.2 Interest in Tandem Repeats

In this section, we summarize theoretical, technical, and medical interests in tandem repeats.

Theoretical Interests.

The abundance of tandem repeats rise some theoretical questions concerning their role in the structure and evolution of the genome. How and why do they appear and evolve? Are they correlated to other local characteristics of the DNA? How frequently do new genes appear through tandem amplification? Already in the 70's, Ohno [Ohn70] argued that gene duplication is a major force in the evolution of genomes. For more information on these topics, the reader may refer to textbooks on molecular evolution like [PH98,Li97].

Technical Interests.

Tandem repeats, especially polymorphic micro- and mini-satellites, have proven useful in many areas of molecular biology. Polymorphic markers are used since the beginning of the 90's to construct low resolution genetic maps. A well-known example is the first genetic map of the human genome built with more than 5000 microsatellites markers [CCW93]. These microsatellites also serve in linkage analysis and positional cloning to detect and locate molecular variations causing disorders [Len02][Chap. 3]. Linkage analysis looks for inheritance correlations between a trait and genetic markers within a pedigree. Polymorphic tandem repeats are markers of choice for Mendelian diseases because the discriminative power of linkage analysis increases with the number of alleles.

In population genetics, polymorphic markers enable biologists to trace the propagation of genetic traits in populations. For instance, highly polymorphic mini-satellites allow to confirm the "Out of Africa" hypothesis, i.e., that our species originated in Africa and invaded afterwards the rest of the world [AAM⁺96]. Differences between alleles of highly polymorphic markers, like the minisatellite MSY1 on the human Y chromosome (see Section 3), give us access to recent populations history.

Because of their level of variability, some polymorphic tandem repeats distinguish any two individuals from the same population and enable the technique of DNA fingerprinting [JWT85a,JWT85b]. Such markers serve as genetic identifiers in forensic studies for the identification of dead corpse, in paternity testing, and so on [GJW85,HJ85]. In 1992, the skeletal remains exhumed in 1985 in Brazil were identified through testing of bone DNA to be those of Dr Josef Mengele, the Auschwitz 'Angel of Death' [JAHS92].

Medical Interests.

At last, tandem repeats are involved in several diseases. Variable minisatellites are known to influence the development of type-1 diabetes, epilepsy and some cancers [BJ97]. Some microsatellites are known to play a role in the regulation of some genes. The most well-known examples are the dozen of severe neurodegenerative diseases caused by large amplifications of CAG/CGG microsatellites either inside or near a gene: fragile X mental retardation, myotonic dystrophy, Huntington's disease, etc (see references in [Wel96,Li97,HGH98]). In healthy individuals, the tandem repeat size varies around a few tens of copies, while in affected individuals the number of copies at the same locus reaches hundreds or a thousand in some cases.

For all these reasons, there are some needs to understand the evolution of tandem repeats. The two problems surveyed in this article should help to establish which mechanisms are responsible for amplification or contraction and in which cases, to estimate how fast the copies of a repeat change over time, to investigate hypotheses on the disease development or to recover recent evolutionary relationships.

1.3 Two problems of interest

Let us first introduce a notation for strings. Let Σ be a finite alphabet of size σ . A sequence of n letters of Σ indexed from 1 to n is called a *word* or a *string* of length n over Σ . We denote the *length* of a word $U := U[1] \dots U[n]$ by $|U|$. For any $1 \leq i \leq j \leq n$, $U[i, j] := U[i] \dots U[j]$ is called a *substring* of U . For $U, V \in \Sigma^*$, $U.V$ denotes the concatenation of U and V . For any integer $h > 0$, U^h denotes the h -th power of U , i.e., the concatenation of h times U . We denote by Σ^* , respectively by Σ^n , the set of all finite words, resp. of all words of length n , over Σ . d_L, d_H denote respectively the Levenshtein and the Hamming distance on Σ^* .

Definitions.

We first define the events a sequence can undergo. The classical point mutation events are **substitution** of a symbol by another, **insertion** or **deletion** of a symbol. Let k be the pattern length, i.e., the minimum size of substrings that can be copied. We term **amplification** the general event that generates copies in tandem of a substring and **contraction** the dual event. The **order** of an amplification is the number of patterns that are copied at a time and its **arity** is the number of copies produced by the amplification plus one.

Definition 1 (Amplification - Contraction). Let T be a text over Σ and $k, i, m > 0$ be integers such that $ik \leq |T|$ and $m \geq 2$. An **amplification of order i and arity m** on T replaces a substring u in T of length ik by u^m of length mik . In other words if $u := T[j, j + ik - 1]$, it creates a new text $T' := T[1, j - 1].T[j, j + ik - 1]^m.T[j + ik, |T|]$. We say u is the pattern of the amplification. A **contraction of order i and arity m** on T' is the dual event of the amplification, that is, it replaces u^m by u and yields a new text $T'' := T$.

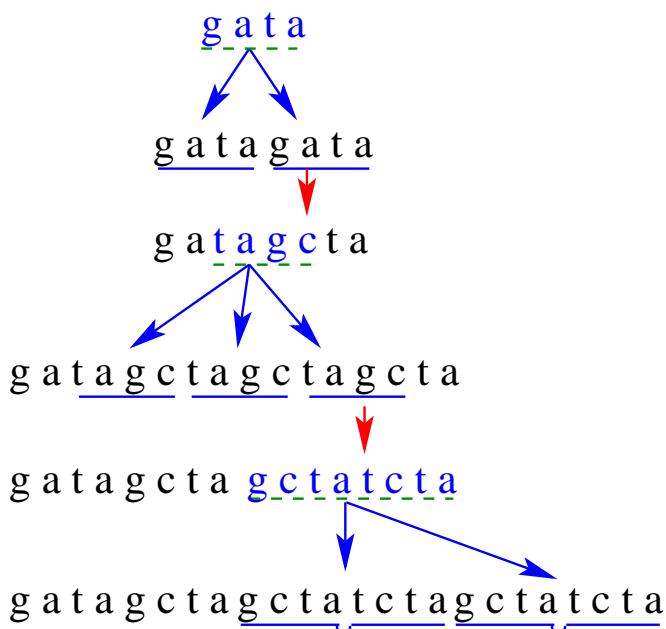


Fig. 2. An example history of an approximate tandem repeat with a basic motif size of $k := 4$. The following events occurred in order: an amplification of order 1 and arity 2 (pattern is *gata*), a substitution of an *a* by a *c*, an amplification of order 1 and arity 3 (pattern is *tagc*), a substitution of a *g* by a *t*, an amplification of order 2 and arity 2 (pattern is *gctatcta*). The patterns and the copies produced are respectively underlined with dashed and straight lines; the number of arrows of an amplification equals its arity.

When amplifications are of arity 2 only, we use the word **duplication** instead. Point mutations and amplifications are illustrated in Figure 2. It is natural to extend this definition by allowing i and m to be rationals instead of integers. Unfortunately, none of the works reported in literature consider this case. Often the terms duplication, triplication, and m -duplication are used instead of amplification. We choose amplification to avoid ambiguity. Also in [BD99], contraction is used with another meaning that we will give later on. In the biological literature, long amplifications are also called expansions.

Given these basic definitions, we can informally state our two problems:

Tandem Repeat History : Given the sequence of one approximate tandem repeat and its minimal pattern size, recover its history of amplifications and mutations.

Tandem Repeat Allele Alignment : Given two allele sequences of the same polymorphic tandem repeat locus, compute an optimal alignment between the two sequences considering point mutations, as well as amplification and contraction events.

2 History of a Tandem Repeat

Given the sequence of an approximate tandem repeat and a minimal pattern size, we want to recover the series of events that led to the present sequence. This problem resembles the one of PHYLOGENETIC RECONSTRUCTION. For this problem, given a set of sequences of the same gene or protein from different species, one wants to recover the evolutionary tree that led to the apparition of the actual species from an ancestral one by a series of speciations (i.e., division of a species in two). The present sequences are associated with leaves and ancestral species with internal nodes of the tree. In our setup, amplifications replace speciations and can have an arity m greater than two (thereby creating m -ary branching). A natural question arises: can one represent the history of a tandem repeat by a tree? It is in general not the case since subsequent amplifications/contractions can act on sequence segments that are not in phase according to the minimal pattern size (cf. Figure 2). In other words for a given pattern size k , all amplifications may not start at position $1 + jk$ for some integer j . To restrict to histories that can be depicted by a tree, one has to limit the starting positions of amplifications to those that respect the pattern phase. This is the **Fixed Boundary** constraint [BD99]. Most researchers envisaged the history problem, which was first proposed by Fitch [Fit77], with this restriction, since they consider the case of tandemly repeated genes where it seems to apply.

Logically, researchers harbored their formalizations from the field of phylogeny reconstruction and considered the problem as an optimization problem with two different criteria. The first criterion is the **maximum of parsimony**; one searches for a tree with sequences labeling internal nodes that minimize the number of evolutionary events. The second is **minimum evolution** and gives rise to distance-based approaches which search for a tree that minimizes the distance between leaves, but do not compute ancestral sequences. [Fit77,BD99,TWY02,EGL02,JKHM02] followed the first line and [TWY02,EG02,EG03] gave algorithms for the second. Note that when one considers tree-like histories as in phylogeny, the order of events is only partially known. Without the help of a molecular clock, it is impossible to order events that lie in different branches of the tree.

All researches achieved on this subject assume that the minimum pattern size, denoted k , is known and that the copies of the repeat have been aligned in a multiple alignment (see Figure 1 for an example). In the case of the distance based approach, the multiple alignment serves to compute a distance between any pair of copies; the method takes as input the resulting distance matrix. This assumption is real restriction since in a complex history there might be several patterns whose sizes are different and not necessarily multiples of another. The problem of multiply aligning sequences is in general NP-hard and thus, the history reconstruction relies on an approximate solution. Finally, this requirement introduces circularity since the multiple alignment itself and the determination of the pattern size depend on the history of the repeat.

The last restriction imposed by all attempts made so far is that apart from point-mutations only amplifications, but no contractions, occurred in the development of a tandem repeat. The history depicts an always increasing repeat which gains new copies at each amplification. From a biological point of view, it is known that tandem repeats undergo multiple amplifications and contractions, but this restriction allows not to consider infinite sequences of events. However for some loci, the reconstructed history seems robust to deletions in the sense that after the deletion of one copy, it still satisfies the constraint of a tandem repeat history. Robustness is also proven when only amplification of order 1 occurred in the history.

After commenting on the different assumptions, we can state the problem more formally. Unfortunately in the literature, most authors do not tackle with exactly the same problem. Often the formalization is expressed in different ways and each work investigates specific versions of the problem. Here, we introduce an unambiguous terminology and a unifying definition such that each version of the problem based on the parsimony criterion is an instance of this definition.

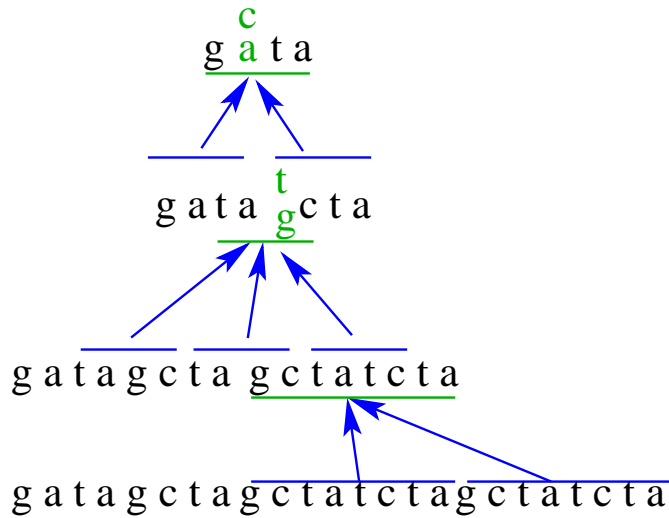


Fig. 3. Reconstruction of a tandem repeat history by successive reductions. The basic motif size is $k := 4$. The history contains three reductions, each combines contractions and mutations if the reduced copies are not identical. The merged copy resulting from the second reduction has a subset of symbols at the first position, while the one of the last reduction has a subset of symbols at the second position. The first ambiguity is resolved in the next reduction and does not appear anymore at the root.

From now on, let k denote the minimum pattern size, T be the approximate tandem repeat sequence, and M be the multiple alignment of the copies of T . Let d be a metric on strings over Σ . Figure 1 shows such a multiple alignment with $n := 8$ copies and $k := 34$ columns. We want to find backwards in time the series of events that led to T from a single copy of an unknown pattern. The history is the repetition of the following process: new identical and adjacent copies were added by amplification and then diverged by point mutations. Recovering the history backwards requires to reduce the number of copies to one. We introduce the notion of **reduction**. A reduction rewinds the process described above. Indeed, it chooses adjacent but (in general) not identical copies, rewinds the differences between copies by point mutations, and contracts the now identical copies. To a **reduction of order i and arity m** , denoted $r_{i,m}$, that transforms $u_1 \dots u_m$ into u is associated a cost function of the form: $C(r_{i,m}) := \sum_j d(u_j, u) + \mathcal{A}(i, m)$ where $\mathcal{A}(i, m)$ is a cost function for an amplification of order i and arity m . The choice of an additive function could be discussed, but seems natural since the genetic events do not occur at the same time and are based on different and independent molecular mechanisms. The reconstruction of a history by successive reductions is illustrated in Figure 3.

Definition 2 (Maximum Parsimony Tandem Repeat History (MP-TRH)). *Let T be a tandem repeat sequence containing n approximate copies of a motif, k be an integer, M be a multiple alignment of k columns and n lines, such that the i -th line contains the i -th copy of T with possibly some indels, and C be a cost function for reductions. The MAXIMUM PARSIMONY TANDEM REPEAT HISTORY problem is to find the minimum cost series of reductions that convert M into a single copy of length k . (The multiple alignment is reduced until one line is left.)*

This series of reductions gives a putative reverse history for T and an associated ancestral pattern. k is also called the repeat unit length and the size of the problem is kn . We denote the n copies of the tandem repeat s_1, \dots, s_n .

Variations Around a Theme.

Some works [BD99,JKHM02] consider an alignment with insertions and deletions. This means that in each line k is the number of columns but not the copy's length. Other authors consider that copies differ only by mismatches/substitutions; in their case, k is the number of columns in the alignment as well as the copy's length. In the first and second cases respectively, the authors use the Hamming and

the Levenshtein distances on Σ^* . When one considers the Hamming distance, one usually discards all columns of the alignment that contain indels. In the case of divergent genes, the number of columns left, k , may be much less than the lengths of the genes that were amplified. In that case, k is not the biological pattern size.

Even if Tang and coworkers [TWY02] include a general cost function in their formalization, all algorithms published so far including theirs only account for differences between copies, but not for amplifications (i.e., $\mathcal{A}(i, m) := 0$).

A major constraint is the **Fixed Boundary** constraint which specifies that amplifications start only in the first column of the alignment. From the computational view-point, it constrains strongly the problem, since many alternative contractions starting at different positions in T , i.e., columns in M , would yield the same result. Also, when boundaries are not fixed, the history cannot be represented by a tree; one needs a more complex structure. From the biological point of view, many tandem repeats do have not an integer, but rather a truly rational number of copies, showing that boundaries of amplifications vary. Nevertheless, when considering tandemly repeated genes, the currently accepted biological model enforces fixed boundaries for amplifications.

Definition 3 (Fixed Boundary Maximum Parsimony Tandem Repeat History (FBMP-TRH)). *This problem is identical to MP-TRH except that reductions start at a position $1 + jk$ for some $0 \leq j \leq n - 2$.*

In each section, we will precise exactly which version of the problem is examined. Under the criterion of Minimum Evolution, the problem is called MINIMUM EVOLUTION TANDEM REPEAT HISTORY and is defined in Section 2.6, where two greedy and one exact algorithm are presented. Its complexity class is unknown. In summary, in the Maximum Parsimony case we face a NP-hard problem for which a 2-approximation, a Polynomial Time Approximation Scheme (PTAS), and several greedy algorithms have been described.

2.1 Benson and Dong’s Approximation and Greedy Algorithms.

[BD99] is the only reference where the general MP-TRH problem with variable boundaries is investigated; the paper also deals with the restricted version with fixed boundaries and amplifications of arity 2 (these are termed “binary duplications” in the paper). Their cost function accounts for the Levenshtein distance between copies and charges no cost for amplifications; that is $d := d_L$ and for any order i and arity m , $\mathcal{A}(i, m) := 0$. The authors provide a greedy algorithm for the general problem, a 2-approximation, and two ways of computing lower bounds. They also report that the greedy algorithm performs in general better than the approximation and is close to the lower bounds. In their paper, the notion of contraction means a reduction.

We describe the greedy algorithm and the 2-approximation which is based on an ordered spanning tree built on the tandem repeat copies. Before that, we give a little correction of the cost function and of the reduction rule².

When a reduction of arity m is applied to the multiple alignment M , a merged copy replaces the m reduced copies. At a column j , differences between the contracted copies may suggest several characters for j -th position of the merged copy. In a reduction step, Benson and Dong authorize such a position to store a set of putative ancestral symbols. To compute optimal ancestral characters, one needs to count for each possible symbol its number of occurrences in the j -th column of the contracted copies and to store the set of all symbols having the majority at the j -th position of the merged copy. If one then chooses any symbol having the majority, say x , for the j -th position of the merged copy, the number of point mutations that must be accounted for is given by: $m - \text{count}[x]$ where the vector count stores the counts mentioned above. Indeed, each contracted copy whose j -th character differs from x requires a single point mutation. The cost function in [BD99] does not consider the relative counts of the characters at a given position; this leads to an incorrect number of mutations, i.e., an incorrect cost.

² G. Benson told us he also noticed this error and corrected it in his algorithms

A Greedy Algorithm for MP-TRH. The greedy algorithm iteratively applies the reduction with the lowest cost ratio; it prefers the reduction with highest arity and breaks other ties arbitrarily. The cost ratio of a reduction is defined as: $\frac{C(r_{i,m})}{k(m-1)}$, i.e., the reduction cost over the arity minus one times the unit length k .

When the arity of reductions is restricted to 2, the greedy algorithm takes $O(kn^3)$ time and when any arity is allowed it takes $O(kn^3 \log(n))$. Let us first explain the complexity for the restricted case. The cost for each possible reduction is computed for increasing order between 1 and $\lfloor \frac{n}{2} \rfloor$. For order i , there are reductions starting at all columns on line j such that $1 \leq j \leq (n - 2ik)$ and at column one on line $(n - 2ik) + 1$. In the left to right order, the cost of reductions whose substrings overlap by all but the first and last characters can be deduced in $O(1)$. Thus, all costs can be obtained in $O(kn)$ time. For each reduction step, the computation takes $O(kn^2)$ and as there are at most $n - 1$ steps the complexity is $O(kn^3)$. In the case of unrestricted arity, the algorithm is more complex since for a given order i and a given position, the costs for arity m can be deduced from those of arity $m - 1$ in $O(1)$. At most n/m costs are updated; the complexity for each step and the total complexity become $O(kn^2 \log(n))$ and $O(kn^3 \log(n))$, respectively.

A 2-approximation for Fixed Boundaries Order-1 Amplifications. The approximation presented in [BD99] is for what the authors call the “restricted problem”: when amplifications have a fixed boundary and their order is 1. In the restricted case, the history can be represented by a leaf and edge ordered and labeled tree whose leaves represent the copies of the tandem repeats. Such a tree is called a **duplication tree**. Internal nodes are labelled by ancestral copies. To an edge (t_j, t_l) , where $t_j, t_l \in \Sigma^k$ are the labels of the linked nodes, is associated the cost $d(t_j, t_l)$.

The proof relies on the concept of **ordered spanning tree** (OST). Given an ordered set of nodes $V := \{1, \dots, n\}$, an OST is a spanning tree on V such that for any two edges (b_1, e_1) and (b_2, e_2) , $b_1 < e_1$, $b_2 < e_2$, we have $(b_1 - b_2)(b_1 - e_2)(e_1 - e_2)(e_1 - b_2) \geq 0$. In other words, if the nodes are placed on a line, the edges are on the same side of the line on the plan and do not cross each other. An example of an OST is given in Figure 4. We will consider OSTs on the ordered set of copies, s_1, \dots, s_n and assign $d(s_j, s_l)$ as the distance between nodes j and l , $1 \leq j, l \leq n$. For both types of trees, the cost of a tree is the sum of the costs of its edges.

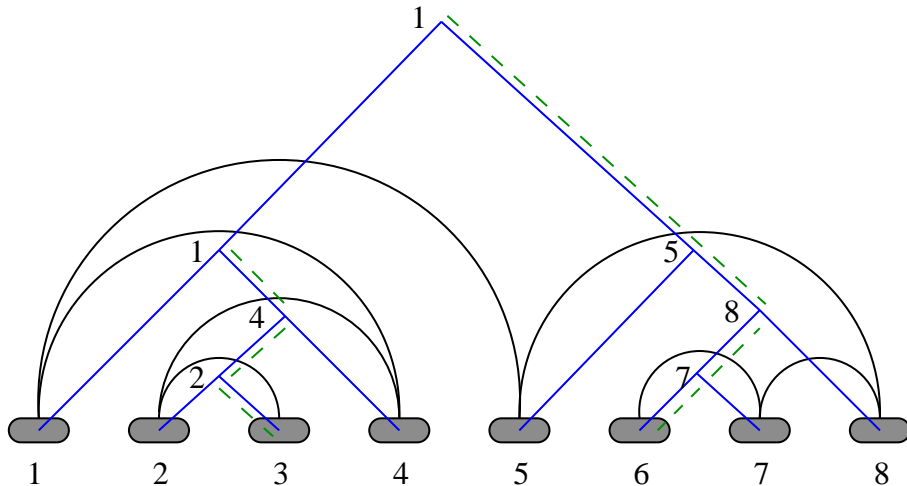


Fig. 4. An ordered spanning tree on $n := 8$ nodes (in curved lines) and its associated duplication tree drawn (in straight lines). The number at an internal nodes represents the leaf label that has been lifted up to that node. In the lifted duplication tree, only edges mark by dashed lines have non-zero cost; all other edges have the same label at each extremity and cost zero.

An ordered spanning tree B on s_1, \dots, s_n describes the structure of a unique duplication tree. It is possible to construct the latter from the former in linear time. This is illustrated in Figure 4. Note that branching the 5th node either to the right or to the left is determined according to the cost of branch entering that node. If in this tree, one assigns to each internal node the label of one of the leaves in its subtree, one obtains a duplication tree B^* whose cost equals the one of the ordered spanning tree. We propose to call B^* a **lifted duplication tree**. The duplication tree given in Figure 4 with the leaf labels associated to internal nodes is a lifted duplication tree. The same idea of lifting the leaf labels up the tree was proposed by Wang and coworkers [WJL96] to compute a 2-approximation for the PHYLOGENETIC MULTIPLE ALIGNMENT problem and their approximate tree was termed a **lifted tree**.

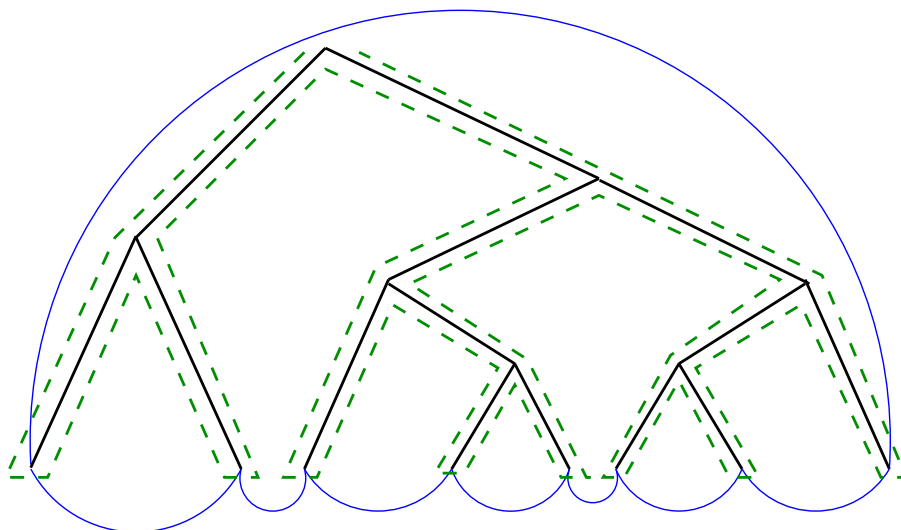


Fig. 5. An optimal duplication tree on $n := 8$ leaves (straight lines), a depth-first traversal of it (dashed lines), and a simple cycle on the leaves (curved lines).

Proof (of the 2-approximation). Let P be the optimal duplication tree, and R be a depth-first traversal of P . As each edge of P is visited twice in R we have $cost(R) = 2 \times cost(P)$. Let SR be the cycle obtained by visiting the leaves in order and cycling from leaf s_n to leaf s_1 . As d satisfies the triangle inequality, $cost(SR) \leq cost(R)$. Now, removing any edge in SR yields an ordered spanning tree TR . Thus, we have $cost(TR^*) < cost(SR) \leq cost(R) = 2 \times cost(P)$; TR^* is a 2-approximation of P . \square

An example of an optimal duplication tree with $n := 8$ copies, its depth-first traversal and a trivial cycle (i.e., corresponding to P , R , and SR , respectively) are given in Figure 5. Note that the topology of TR^* varies very little with the input (since TR equals SR minus one edge) and presents no biological interest in itself. On the opposite, a true minimum OST and its associated duplication tree provide an informative approximate solution. Benson and Dong exhibit a dynamic programming recurrence to compute the minimum OST (on all possible intervals of lines in M) in $O(kn^2 + n^3)$. We will see below that Tang et al. [TWY02] also use the lifting technique to obtain a 2-approximation for this problem in $O(n^2(k + n^3))$.

Open questions:

- Is MP-TRH with variable boundaries approximable?

2.2 Tang et al. Dynamic Programming Approach.

In their framework, Tang and coworkers [TWY02] allow only duplications, that is amplifications of variable order but of arity 2, and fixed boundaries. They generalize the duplication tree introduced by [BD99] in a **duplication model** to enclose information relative to duplications of order higher than one. We first describe this model. We present their dynamic programming scheme for what they call the SINGLE GENE DUPLICATION problem, that is FBMP-TRH restricted to arity 2 and order 1. (Actually, the authors do not really tackle with their general model.) In Benson and Dong’s vocabulary, this corresponds to the “restricted” FBMP-TRH. The dynamic programming scheme is flexible and by combining it with the lifting technique, they obtain a 2-approximation and claim a PTAS can be achieved with the same technique. The authors also give an algorithm for the distance-based problem, but we will delay its presentation until Section 2.6.

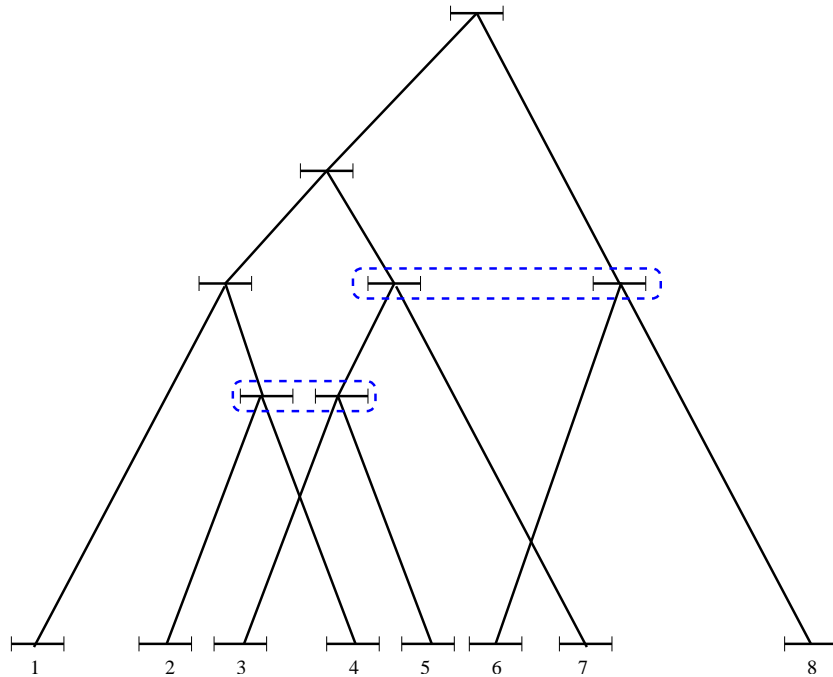


Fig. 6. An example of a duplication model or of an unrooted duplication tree with $n := 8$ leaves. Each node represents a copy of the tandem repeat, either an ancestral one at an internal node or a present one at a leaf. Subsets of internal nodes that are duplicated together by a duplication of order 2 are surrounded by a dashed line. These correspond to blocks in the terminology of [TWY02].

The Duplication Model. The duplication model is a duplication tree (from [BD99]) with higher order duplications. A duplication of order i duplicates i adjacent copies, say $s_1 \dots s_i$, into $s_{1,l} \dots s_{i,l} s_{1,r} \dots s_{i,r}$ ³ where $s_{j,l}, s_{j,r}$ denote the right and left children of s_j . As a duplication tree is ordered, the left to right order of nodes represents the sequence order of the copies; this is valid for both internal nodes and leaves. If we place a duplication of order i in such a tree, it must take i nodes representing adjacent copies and create i left and i right children. Therefore for $j < h \leq i$, the edge connecting s_j and its right child $s_{j,r}$ crosses all edges connecting s_h and its left child $s_{h,l}$, i.e., $(s_h, s_{h,l})$. This is the only way two edges can cross in the model and it requires that $i \geq 2$. Such a model is said to be **consistent** with the ordering of the leaves/copies on the sequence. This is the main constraint that a tree must satisfy

³ Note that there are no commas in the notation; it means that the copies are adjacent on the chromosome.

to be a duplication model. Nodes involved in higher order duplications form a **block** and all blocks are memorized in the duplication model. Figure 6 illustrates the notion of duplication model.

For a given duplication model, it is possible change the order of nodes to uncross all edges and obtain a unique planar tree topology. Tang et al. name it the associated **phylogeny**. The term phylogeny means a binary tree without order constraint that represents the evolution of genes/species associated with its leaves. A way to answer our problem is to use a phylogenetic reconstruction method on the set of copies without the order constraint, and check if the returned phylogeny has an associated duplication model. There exist efficient heuristics for phylogenetic reconstruction (for instance [Gas97]). Tang et al. investigate what we call the **Recognition problem**, that is to find the duplication model associated with a rooted phylogeny if it exists, and they give a $O(n^2)$ algorithm for it. Note that the problem size is n , i.e., the number of leaves of the phylogeny.

The Dynamic Programming Scheme. In [TWY02], the authors define the MP-TRH with fixed boundaries with an arbitrary additive cost function for the reductions. However, the dynamic programming scheme they present is further restricted to the FBMP-TRH with order 1 duplications only. In this case, the duplication model has no edges crossing each other and is simply a tree on an ordered set of leaves. In the remainder of this section, we use the word “tree” instead of “duplication model” for simplicity. The authors consider the Hamming distance between words, and not the Levenshtein distance as in [BD99].

Let s belong to Σ^k , i, j be integers such that $1 \leq i \leq j \leq n$ and S, S_l, S_r be subsets of Σ^k . A tree for an interval subset of the copies, s_i, \dots, s_j , is said to **span** the interval $[i, j]$. With Maximum Parsimony criterion, each internal node is labelled with a putative ancestral copy. We introduce a notation used below in the recurrence. Let $D([i, j])$ be the cost of an optimal tree spanning $[i, j]$. Let $D([i, j], s)$ be the optimal cost of a tree spanning $[i, j]$ and whose root is labelled by s . If $i < j$, let an integer m satisfy $i \leq m < j$ and $D([i, j], s, m)$ be the minimum cost of tree spanning $[i, j]$ whose root is labelled by s , and whose left and right subtrees span $[i, m]$ and $[m + 1, j]$, resp. For $s \in \Sigma^k$, $D([i, i], s)$ is initialized to 0 if $s = s_i$ and to infinity otherwise. It follows from these definitions that:

$$D([i, j]) = \min_{s \in S} D([i, j], s) \quad (1)$$

$$D([i, j], s) = \min_{i \leq m < j} D([i, j], s, m) \quad (2)$$

$$D([i, j], s, m) = \min_{v \in S_l} (D([i, m], v) + d(s, v)) + \min_{w \in S_r} (D([m + 1, j], w) + d(s, w)) \quad (3)$$

A tree spanning the whole repeat can be recovered by backtracking through the matrices from entry $[1, n]$. If one chooses S, S_l, S_r equal to Σ^k the returned cost and tree are optimal, but the running time is exponential in k . Indeed, it is in $O(|\Sigma|^{2k}(k + n^3))$ since computing all Hamming distances between any two strings of length k is done in $O(k \times |\Sigma|^{2k})$ and filling the matrices for each of the n^2 intervals takes $O(|\Sigma|^{2k}n)$ time.

Now, if one restricts S, S_l and S_r in such a way that to an internal node is associated the label of either its left or of its right child, the result is a lifted duplication tree. By a lemma from [WJL96], the cost of such tree is at most twice the cost of an optimal tree. This gives a 2-approximation algorithm in $O(n^2(k + n^3))$ time. The authors also report that the PTAS developed by Wang and Gusfield [WG97] can be adapted to this problem.

The 2-approximation of [TWY02] yields a lifted tree as the 2-approximation of [BD99] obtained from the trivial cycle. The latter is more effective since it requires linear time. Its improvement using the minimum OST is computed in only $O(kn^2 + n^3)$ time. Another remark is that the dynamic programming scheme delivers an optimal solution in a time that is only exponential in k . Thus, FBMP-TRH with order 1 duplication is Fixed Parameter Tractable for parameter k . We refer to [DF99] for details on Parameterized Complexity.

2.3 An Exhaustive Exploration Approach

Elemento and coworkers [EGL01,EGL02] investigate the same problem as in [TWY02]: recovering the duplication tree of tandemly repeated genes. They define independently the concept of duplication model and call it a **partially ordered duplication history**. They name such an unrooted history a **tandem duplication tree**. The authors argue that the Fixed Boundaries restriction applies here, because the main mechanism of gene amplification is unequal recombination. They exhibit an exponential algorithm, named *DTExplore*, that exhaustively searches the space of duplication histories. The algorithm works for a limited number of copies, in the order of $n = 10$. In addition to some applications, the algorithm is used to compare empirically the number of phylogenies with n leaves with the number of duplication models with n copies. The authors conclude that, although the number of duplication models seems exponential, it represents a small fraction of the number of phylogenies. This question is addressed later on in [GHJMM03] and is detailed in Section 2.4. We also present their $O(n^2)$ solution for the RECOGNITION problem as defined in Section 2.2.

Like in [TWY02], authorized amplifications are duplications (i.e., of arity 2) of variable order. In their vocabulary, a duplication of order x is denoted an x -duplication. The distance measure between gene copies is the Hamming distance and no cost is charged for amplifications (i.e., $d = d_H$ and $\mathcal{A}(i, m) = 0$ for all i, m).

An Exhaustive Search Strategy. The algorithm *DTExplore* simulates the duplication process to explore in a depth-first search manner the space of all duplication histories of n leaves. It starts with a rooted tree with two leaves and applies a duplication of order x to obtain a tree with $2 + x$ leaves. The process is iterated until the number of leaves reaches n . The current phylogeny is given as input to Fitch's algorithm [Fit71] which computes its Maximum Parsimony score (i.e., the optimal ancestral copies for all internal nodes). The algorithm backtracks and duplications of any possible order are applied to visit all topologies. *DTExplore* outputs the phylogenies with the minimal score.

This first version of *DTExplore* suffers from redundancy as it generates several times the same rooted duplication model. An improvement is achieved by enabling *DTExplore* to memorize which duplication model has already been visited. This is performed by encoding each model in a character string and storing all codes in a prefix tree. Each generated topology is first encoded, searched for in the prefix tree, and its maximum parsimony is evaluated only if it was not found in this data structure. The speed improvement is drastic.

A Simple Recognition Algorithm. Given the topology of a phylogeny on n leaves, the RECOGNITION problem is to decide if the topology also is a duplication model. Elemento et al. [EGL02] report a simple $O(n^2)$ algorithm for it. It iteratively reduces the phylogeny by replacing the $2i$ leaves of an order i duplication by their fathers. At most n such steps are performed, if all duplications were of order 1. At each step, to identify a possible duplication, it searches for a subset of adjacent $2i$ leaves that are appropriately intermingled. If this search fails, the phylogeny is not a duplication model and the algorithm returns false. Otherwise, it stops with a tree reduced to a root and answers true. The search is done in $O(n)$ time at each step. After investigating the combinatorics of duplication models, Gascuel et al. [GHJMM03] improve the algorithm's complexity to $O(n)$ (cf. Section 2.4).

2.4 Combinatorics of Duplication Trees

In the same framework than [EGL02], Gascuel and coworkers [GHJMM03] investigate the number of different rooted duplication trees with n leaves, denoted $\text{RDT}(n)$. They exhibit a recurrence for it and show it is the double of the number of unrooted duplication trees. They deduce an algorithm to uniformly sample duplication trees and a linear time procedure for the RECOGNITION problem. Here again, amplifications are of variable order and of arity 2.

Counting Duplication Trees. We review the main recurrence for rooted duplication trees. To obtain this recurrence, the authors introduce the notion of an **(l, i) duplication** where i is the order and l the number of copies located after the last copy involved in the duplication. So, an (l, i) duplication duplicates i copies $s_{n-l-i+1} \dots s_{n-l}$, where i and l satisfy $1 \leq i \leq n$, and $0 \leq l \leq n - 2i$. Given a rooted duplication tree R , an (l, i) duplication is said to be **visible** if none of the $2i$ copies it created has been further duplicated in R . Let $P(n, l)$ be the subset of $\text{RDT}(n)$ whose leftmost visible duplication is an (l, i) duplication for some i such that $1 \leq i \leq (n - l)/2$. Let $p(n, l)$ be the cardinality of $P(n, l)$. By definition, for $l > n - 2$, $P(n, l) = \emptyset$, $p(n, l) = 0$ and $p(2, 0) = 1$.

Theorem 1. *Let $n > 2$ and $0 \leq l \leq (n - 2)$. $P(n, l)$ and $\cup_{j=0}^{l+1} P(n - 1, j)$ are in one-to-one correspondence.*

Proof. Let us denote the substring of $2i$ copies created by an (l, i) duplication in the present sequence by $s_{n-l-2i+1} \dots s_{n-l-i} s_{n-l-i+1} \dots s_{n-l}$. For all $1 \leq f \leq i$, $s_{n-l-i+f}$ is the twin of $s_{n-l-2i+f}$, i.e., they are offspring of the same father copy. Let $T \in P(n, l)$; deleting s_{n-l-i} , the left child of the rightmost copy duplicated by the (l, i) duplication, in T maps T to T' . If $i = 1$ then the leftmost visible duplication in T' is at the right of s'_{n-1-l} and thus, T' is in $P(n - 1, j)$ for some j in $0 \leq j \leq l$. If $i > 1$ the (l, i) duplication in T becomes the leftmost visible duplication in T' and is an $(l + 1, i - 1)$ duplication; so T' belongs to $P(n - 1, l + 1)$. As the transformation is reversible the mapping is a bijection. \square

As all trees have a leftmost visible duplication, if we denote the cardinality of $\text{RDT}(n)$ by $\text{rdt}(n)$, we have $\text{rdt}(n) = \sum_{l=0}^{n-2} p(n, l)$. Combined with the recurrence of Theorem 1 it provides a way to compute $\text{rdt}(n)$.

Although the root of a duplication tree is necessarily located on the path between the left- and right-most copies ([Fit77]), not all edges on this path are possible root locations. First, the root cannot be below some duplication of order strictly greater than 1. Second, it comes out that in average only two locations are possible. Hence, the surprising result that the number of unrooted duplication trees with n leaves equals $\text{rdt}(n)/2$. Asymptotically, this cardinality behaves like $(\frac{27}{4})^n$ when n tends towards infinity.

An $O(n)$ Recognition Algorithm. The concept of visible duplication allows to improve the recognition algorithm described in Section 2.3. It proceeds by iteratively agglomerating leaves of a duplication (such leaves belong necessarily to a visible duplication). The improvement consists in choosing the leftmost visible duplication at each step; an amortized analysis shows that all steps take $O(n)$ time altogether. For this, the scan for a duplication proceeds from left to right and the endpoints of already encountered blocks are memorized.

Open questions:

- In the case of fixed boundaries, study the number of amplification trees when amplifications of variable arity are allowed.

2.5 Complexity, Approximability and Other Results

In [JKHM02], Jaitly et al. consider the FBMP-TRH, with a maximum arity 2, with the Levenshtein distance between strings $d := d_L$, and no cost for amplifications. They proved that this restriction of FBMP-TRH is NP-hard by reducing it to the MAX-CUT problem. They exhibit a PTAS that uses the lifting technique to partition the topology and exact optimization by dynamic programming to compute optimal labels for subtrees of constant size. A detailed sketch of the proof is given in [JKHM02] and it relies on previous difficult approximation results for the PHYLOGENETIC MULTIPLE ALIGNMENT problem by Wang and coworkers [WJL96, WG97, WJG01]. As in [BD99], the authors notice the relation of FBMP-TRH with the ORDERED LEAVES STEINER TREE problem. They show their algorithm also is a PTAS for the latter.

Zhang et al. [ZMW02] describes a $O(n)$ algorithm for the RECOGNITION problem from rooted phylogenies (with duplications of variable order, the distance used over Σ^* is not given). Assume the leaves of the phylogeny are numbered in order from 1 to n . Their method identifies for all leaf $j := 1, \dots, n$ the internal nodes belonging to duplications of order ≥ 2 (to blocks in the terminology of Tang et al.) that allow leaves j and $j + 1$ to be put next to each other in the duplication model. For this, they associate to each node v the pair $(l(v), r(v))$ where $l(v), r(v)$ denote resp. the smallest and largest leaf numbers in the subtree of v . Fast identification of block's nodes is achieved by comparing pairs $(l(v), r(v))$ of nodes on the path to leaf j to those of nodes at the same level on the path to leaf $j + 1$. Given an unrooted phylogeny, the search for a duplication model has to be performed for all $O(n)$ possible root locations because in practice one infers unrooted phylogeny. Thus, although in linear time, their algorithm is less effective than the one of Gascuel et al. [GHJMM03] that deals with both case of rooted and unrooted phylogenies.

Zhang and coworkers also proposed a greedy search strategy for the FBMP-TRH. It infers a phylogeny with a traditional reconstruction method, checks if it is associated to a duplication model, and if not attempts a transformation of the topology (like a Nearest Neighbor Interchange). The two last steps are iterated until a duplication model is found. Such a strategy was first developed in the field of phylogeny reconstruction. The authors report it performs better than Benson and Dong's greedy algorithm or Tang et al.'s Window method on three real data cases.

Open questions:

- Does the MP-TRH with variable boundaries admit a PTAS?
- For some parameters, are these problems fixed parameter tractable in the sense of Fellows and Downey's theory of parameterized complexity [DF99]?
- Improve the Zhang et al.'s greedy strategy by inferring directly a duplication model and invent a transformation operation that allows to visit each possible duplication model on n leaves.

2.6 Minimum Evolution Approaches

As mentioned above, the criterion of Minimum Evolution borrowed from Phylogeny leads to a formalization of the problem that differs from MP-TRH. The data is a matrix D giving the pairwise distance between any pair of copies in T and the ordered list of copy numbers. The output is a duplication tree whose sum of the branch lengths is minimum. We call this problem the MINIMUM EVOLUTION TANDEM REPEAT HISTORY.

This formalization implicitly considers fixed boundaries since the history is represented by a tree. It is thus not as general as MP-TRH. Moreover, as the input is a distance matrix, the sequences are disregarded in the remaining of the algorithms. The link with combinatorics and algorithmics on words resides only in the computation of the input pairwise distances. Nevertheless, we include a section on these methods for the sake of completeness and because such approaches were shown to be reliable in practice [EG02].

The two methods presented here, *DTScore* from [EG02] and *Window* from [TWY02], optimize a local criterion on the tree and are based on the same algorithmic scheme. A current list of leaves in the tree is maintained and initialized to the original list of copies. The methods iterate an agglomeration step that replaces a subset of $2i$ adjacent leaves resulting from an order i duplication by new leaves representing their parent copies. Such a subset of $2i$ leaves is called a **window** in [TWY02]. Entries of the distance matrix corresponding to the deleted copies are removed and entries for the parent copies are inserted. Entries for the parent copies are the average of their children distances. The algorithms proceed until 2 or 3 leaves are left. The window is chosen to maximize a score function. *DTScore* and *Window* differ in the choice of this function. *Window* considers the average of the Hamming distances between twin copies j and $j + i$ for all possible j . This is known to produce the correct tree if the data respect the **Molecular Clock** hypothesis, i.e., if evolution proceeded at the same pace in each branch of the tree. In practice, it is not often the case. In *DTScore*, the score of a pair $(j, j + i)$ is the number of times j and $j + i$ are next to each other in every possible quartet of leaves $(j, j + i, l, m)$ according to the **Four-Point Condition** ([Bun74]). The score of a window is the minimum of the scores of its

pairs; the window with the maximum score is selected. This scheme works even if the Molecular Clock is not respected.

Empirical tests on pseudo-randomly generated trees and sequences confirm *DTScore*'s improved ability to recover the correct tree compared to *Window* and to maximum parsimony approaches (cf. [EG02]). For both programs, the time complexities is $O(n^4)$. Elemento and Gascuel [EG03] consider the ME-TRH under the global criterion of the ordinary least square errors. They propose an exact algorithm in $O(n^3)$ time and $O(n^2)$ space when the duplication order is restrained to 1.

Open questions:

- Find an algorithm to optimize ordinary least square errors when the order is variable.

2.7 Biological Validations

Most of the literature referenced for the TANDEM REPEAT HISTORY problems includes tests on real data sets, mainly on clusters of tandemly repeated human genes. For instance, Zinc Finger genes in [TWY02,ZMW02], genes for Olfactory Receptors and the internal tandem repeat of a mucin gene MUC5 in [ZMW02], as well as two immunological gene loci, TRGV and IGLC, in [EGL02]. In all cases, the authors found the computed histories were consistent with what is known about the evolution of the gene family and most histories include mainly order 1 duplications.

The most detailed discussion on biological validation is in [EGL02]. The mechanism for gene duplication is hypothesized to be unequal recombination. In the case of the TRGV locus where $n := 9$, the history reconstructed with *DTExplore* is also found by *DNAPENNY*, an exact phylogenetic reconstruction method ([HP82]), although the probability of finding a duplication model when exploring the space of phylogenies is already low for $n = 9$ (see Section 2.4). Moreover, the most recent duplication in the tree has order 2 and corresponds to a polymorphism observed in some human populations. Indeed, the two additional gene copies are missing in individuals of the Tunisian, Lebanese, French, Black-African and Chinese populations. These two evidences combined with a bootstrap test constitute a strong validation of the model with respect to the mechanism of duplication.

With the availability of complete genomes, numerous tandem gene clusters are discovered and represent potential data for such analysis. We believe biologists will investigate tandem duplication history once the methods described here have been advertised more widely.

3 Allele Alignment

Among polymorphic tandem repeats, hypervariable minisatellites cumulate variations in their number of copies, as well as in the sequence of the repeats. They belong to the most polymorphic markers: at a single locus, one encounters much more alleles than for other loci. Hence, the numerous differences between alleles provide us with detailed information on the evolutionary processes at work. To understand variation in the sequences of a tandem repeat locus, we need to be able at least to compare alleles in a pairwise manner. If one can measure the dissimilarity by a metric, it becomes possible to infer evolutionary relationships of a set of alleles using distance-based phylogenetic methods. Even better would be to simultaneously align several alleles and then use Maximum Parsimony phylogenetic reconstruction to compute a tree with ancestral alleles and to count mutations along the branches. In this section, we consider the problem of PAIRWISE ALLELE ALIGNMENT. At present, MULTIPLE ALLELE ALIGNMENT remains an interesting future work.

Bérard and Rivals [BR02,BR03] introduced this problem and noted that the events are not commutative, which introduces a major difficulty compared to classical SEQUENCE ALIGNMENT [SK99]. In the case of amplifications and contractions of order 1 and arity 2, they describe an exact algorithm in $O(\max(m, n)^4)$ time, where n, m denote the sequences lengths. The alignment distance is a metric and can serve as input for distance-based phylogenetic reconstruction. In [BR03], the method is applied to a set of alleles from the MSY1 locus, a minisatellite on the human Y chromosome. These biological

experiments demonstrate that the approach enables recovering allele relationships. It is shown for instance that phylogenetic alleles having a common origin are segregated according to the population of their bearer. The approach used to compare minisatellite alleles is also valid for other types of markers.

In the remaining of this Section, we present the notion of minisatellite maps and the hypothesized evolutionary model, we formally define the problem, and show that the non-commutativity forces us to consider the order in which events occur to find an optimal alignment. We then describe the solution that combines dynamic programming and computations of maximum independent sets in overlap graphs. We conclude with biological validations.

3.1 Minisatellites Maps and the Evolutionary Model

Along the tandem array of a minisatellite, the repeated unit varies in sequence. For a fixed phase, the different adjacent substrings are called **variants** of the repeated unit. In 1991, Jeffreys and colleagues design a reaction to obtain the sequence of variants of the array. This specific Polymerase Chain Reaction (PCR) is called the **Minisatellite Variant Repeat** reaction or MVR-PCR for short [JMT⁺91]. It yields a sequence in which each variant/substring is encoded by a symbol. Such sequences are called **minisatellite maps**. This technology permitted investigations of minisatellites instability processes. Here we consider that alignment is performed on allele maps and not on their DNA sequence. From the computational view-point, maps are still sequences over a finite alphabet, which is not the DNA alphabet. In the sequel, we use the terms string, sequence or map without distinction. It also means that boundaries of the variants are fixed by this technology.

The evolutionary model defines the set of events a map can undergo and associate a real cost with each event. Here, we consider point-mutations, amplifications and contractions of order 1 and arity 2 as defined in Section 1.3. Note that events operate on variants instead of on DNA bases; e.g., a mutation changes a variant into another, a deletion removes a complete variant from the map, etc. To each event is associated a fixed real cost that we denote: A, C, D, I, M for resp. amplification, contraction, deletion, insertion and mutation. We assume our model is symmetrical: dual events have equal costs ($A = C, D = I$). To stick to biological conditions, we assume $A, C < M, D, I$. Note that if $I > A + M$ then an amplification followed by a mutation is always preferred to an insertion. Without loss of generality, we make this hypothesis on costs.

Individual 1		Individual 2	
Event	Sequence	Event	Sequence
	a a a a a		a a a a a
mutation	a e a a a	mutation	a a a b a
		5*amplification	a a a b b b b b a
		mutation	a a a b b c b b b a
		mutation	a a a b b c b d b a
		amplification	a a a b b c b d d b a

Fig. 7. Example of evolution of a minisatellite in two individuals.

3.2 A Dynamic Programming Approach Combined With Graph Algorithms

From now on, let r, s be two maps resp. of length n and m over the alphabet Σ . Let $-$ be an additional symbol that does not belong to Σ and will denote an inserted or deleted position in an alignment. An edit script between r and s is a sequence of events that transform r into s . An alignment of r and s is a representation of an edit script that respects the order of the positions in r and s . The alignment cost is the sum of its operation costs where aligning two identical variants costs zero. It is shown that the alignment cost is a distance metric [BR03]. An example of evolution of different alleles from the same

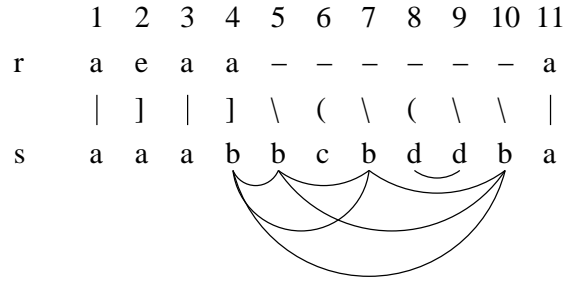


Fig. 8. Alignment of the maps $r := aeaaa$ and $s := aaabbcbddb$. The arch $bbcbddb$ and its inner-arches are drawn by curved lines under s . In the middle line, ‘|’, ‘]’, ‘\’, ‘(’ denote resp. a match, a mismatch, an amplification, and an amplification+mutation.

ancestor allele is given in Figure 7 and an optimal alignment for the two resulting maps is shown in Figure 8.

Definition 4 (Allele Alignment). *Given two maps r and s respectively of length n, m over an alphabet Σ and a alignment scoring scheme, find the alignment between r and s that has minimum cost.*

In Figure 7, the substring of s from position 4 to 10, $bbcbddb$, shows that the final variant at each position does not appear in the order of the sequence: at some stage, position 8 has still the ancestor state b and not its final state d , while position 10 is already a b . The minimal cost series of events to create such a substring is order-dependent. If one computes incrementally alignment for longer and longer prefixes, we cannot find the optimal order of events. This happens when aligning a position in r with a substring of s having identical first and last characters. Such a substring is called an **arch**. The first position of the arch in s is aligned to the position in r , and all other arch’s characters are generated in s from the first position. The authors prove that at least in an optimal **generation** of the arch, the last character is obtained by an amplification of the first when the positions are adjacent, and then all other positions in between are generated afterwards. Thus, one needs to consider the arch as whole, and not in order of increasing prefixes. Such an arch generation avoids a mutation per arch. An arch may include other inner-arches, but an alignment cannot include all possible arches. Because of the optimal order of generation, two arches cannot belong to the same alignment if they are **incompatible** with each other, i.e., if they overlap each other by more than one variant (here overlap means that strict inclusions are allowed). It follows that an optimal arch generation should contain the largest number of pairwise compatible arches. The symmetrical situation, when the arch in r is aligned to a single symbol in s , is called an arch **compression**.

Arches represent intervals of a map and incompatibility defines an overlap relationship between these intervals (not an overlap+containment relationship). Consider the graph G whose nodes are arches and whose edges link two nodes if their arches are incompatible. G is an overlap graph. It is shown that computing the maximum subset of compatible arches is equivalent to finding a maximum independent set in G .

A preprocessing procedure will take as input a map t and compute for each possible interval the cost of the corresponding arch generation or compression (if any) using an adaptation of Apostolico and coworkers’s algorithm [AAH92]. This procedure is applied to s and to r and results are stored in two matrices that require quadratic space.

The problem is solved by filling a dynamic programming matrix \mathcal{A} whose entry $\mathcal{A}(i, j)$ is the optimal alignment cost between the prefixes of s and r of length i and j resp. Five dependencies between $\mathcal{A}(i, j)$ and its direct neighbors account for the five possible evolutionary events: amplification, amplification+mutation (A_M), mutation/match, contraction, and mutation+contraction (M_C). (Note that there are no insertion, nor deletion because of our hypothesis on costs.) Moreover, up to i , resp. j , dependencies with non adjacent entries on the same line, resp. on the same column, account for arch generations, resp. compressions. These costs are denoted G_l and K_l . These dependencies can be

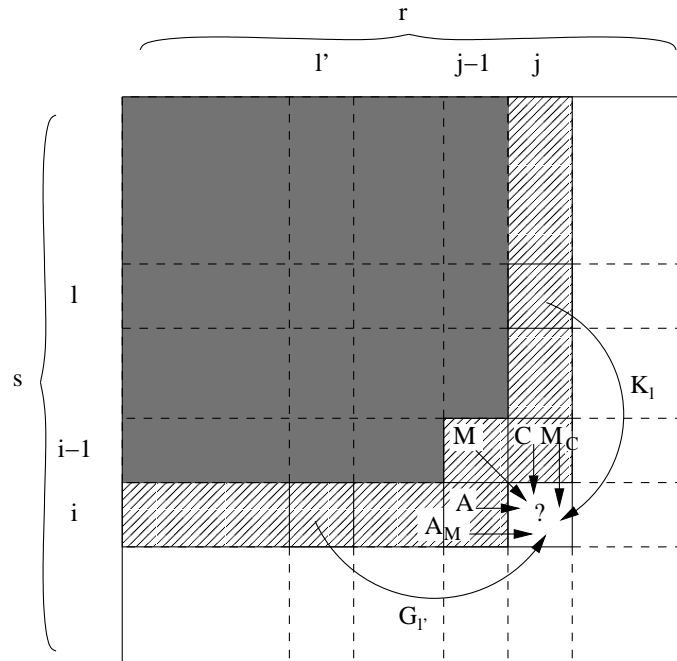


Fig. 9. Dependencies in the dynamic programming matrix. To compute cell $\mathcal{A}(i, j)$, we need at most all cells in the striped patch but not the ones in the dark patch. Dependencies are shown by arrows. For arch generations and compression the arrows are indexed by the beginning position of the arch to show that there are multiple dependencies.

evaluated in $O(1)$ time thanks to the preprocessing. The dynamic programming recurrence is illustrated in Figure 9. The algorithm requires $O(\max(m, n)^4)$ time and $O(\max(m, n)^3)$ space.

Open questions:

- What is the complexity class of PAIRWISE ALLELE ALIGNMENT when higher order amplifications and contractions are allowed? Find a practical solution.
- An algorithm for MULTIPLE ALLELE ALIGNMENT under any of the assumptions mentioned above.

3.3 Biological Validations

The algorithm was implemented in a program named *MS_ALIGN* and applied to alleles of the human minisatellite MSY1 [JBT98]. MSY1 is a hypervariable minisatellite locus on the human Y chromosome. Its repeat unit is 25 base pairs long and five different variants, which differ from each other by at most 3 substitutions, have been observed. Amplifications and contractions of order 1 and arity 2 were shown to be the most probable events experimentally [ALO02]. This is in agreement with the model. The data set comprises 609 alleles from all over the world, with the corresponding map, as well as the population of origin and the Y-chromosomal genetic group (the technical term is **haplogroup**) if known. All pairwise distances between alleles were computed. The experiments consisted in reconstructing phylogenetic trees with a distance-based method for all or a subset of the alleles using the corresponding distance matrix. The resulting trees are confronted to other experimental data.

First, an evolutionary tree of the haplogroups was obtained from MSY1 with average distances and found to resemble strongly the tree reconstructed from other less polymorphic markers. However, the MSY1 tree gives a higher resolution. Second, the authors looked at the trees of all alleles from a given haplogroup. In many such trees, alleles are grouped by populations of origin when these are geographically distant. This is consistent with the already observed geographical specificity of the Y chromosome [JBT98]. An example of tree for the fourth haplogroup is given in Figure 10. In there, the Japanese population is perfectly separated from the Tibetan and Mongolian alleles in this haplogroup.

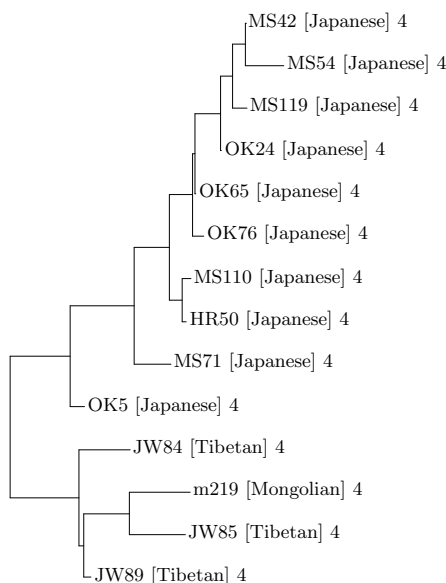


Fig. 10. Phylogenetic tree of haplogroup 4 built from a distance matrix produced by *MS_ALIGN*. Each individual is represented by its code, its population origin and its haplogroup.

4 Conclusion

Organisms have the possibility to locally duplicate, triplicate, etc. a segment of their genome, and also to remove one or more copies among adjacent identical segments. This creates repeats with varying numbers of copy units, called tandem repeats. As point mutations also alter the copy units, tandem repeats display variation in copy sequence and in length. It follows that a tandem repeat has a history and that any two individuals may have different tandem repeat sequences, alleles, at the same genomic location. For many reasons, biologists are interested in tracing back the history of a tandem repeat and to compare different alleles of a tandem repeat. In this paper, we surveyed these two problems, TANDEM REPEAT HISTORY and TANDEM REPEAT ALLELE ALIGNMENT, which are related to the evolution of tandem repeats. We gave an overview of algorithmic and combinatorial results on these topics, as well as detailed biological motivations. We provided a unified framework for their formalization. Among the algorithms mentioned in this survey, some have been implemented and are available on the Internet at <http://degas.lirmm.fr/REPSEQ:DTScore> and *MS_ALIGN*.

Surely, the solutions for TANDEM REPEAT HISTORY would improve in biological significance, if the algorithms could consider both amplifications and contractions. In order to avoid considering infinite histories, a relevant constraint would be to limit the number of contractions: not all copies created by an amplification can later be removed by a contraction. The TANDEM REPEAT ALLELE ALIGNMENT problem is exactly solved when amplifications and contractions are limited to order 1. Relaxing these constraints represents future lines of research. Combinatorial properties of tandem repeat histories also require more investigations and could lead to algorithmic improvements. Although probabilistic approaches based on Maximum Likelihood are recognized as the most reliable methods in phylogeny, they were neglected until now for TANDEM REPEAT HISTORY. The design probabilistic methods for this problem also seems a promising direction of research.

References

- [AAH92] Alberto Apostolico, Mikhail J. Atallah, and Susanne E. Hambrusch. New clique and independent set algorithms for circle graphs. *Discrete Applied Mathematics*, pages 1–24, 1992.

- [AAM⁺96] J. A. Armour, T. Anttinen, C. A. May, E. E. Vega, A. Sajantila, J. R. Kidd, K. K. Kidd, J. Bertranpetit, S. Paabo, and A. J. Jeffreys. Minisatellite diversity supports a recent African origin for modern humans. *Nat Genet*, 13(2):154–60, 1996.
- [ALO02] R. Andreassen, J. Lundsted, and B. Olaisen. Mutation at minisatellite locus DYF155S1: allele length mutation rate is affected by age of progenitor. *Electrophoresis*, 23(15):2377–83, Aug 2002.
- [BD99] Gary Benson and Len Dong. Reconstructing the Duplication History of a Tandem Repeat. In *ISMB*, pages 44–53, Heidelberg, Germany, 1999.
- [Ben99] G. Benson. Tandem Repeats Finder: a Program to Analyze DNA Sequences. *Nucleic Acid Research*, 27(2):573–80, 1999.
- [BJ97] J. Buard and A. J. Jeffreys. Big, bad minisatellites. *Nat Genet*, 15(4):327–8, 1997.
- [BR02] Sèverine Bérard and Eric Rivals. Comparison of Minisatellites. In G. Myers, S. Hannenhalli, S. Istrail, P. Pevzner, and M. Waterman, editors, *Proc. of the Sixth Annual International Conference on Computational Molecular Biology*, pages 67–76, Washington DC, USA, 2002. ACM Press.
- [BR03] Sèverine Bérard and Eric Rivals. Comparison of minisatellites. *J. of Computational Biology*, 10(3-4):357–372, 2003.
- [Bun74] P. Buneman. A note on metric properties of trees. *J. of Combinatorial Theory series A*, 17:48–50, 1974.
- [CCW93] D. Cohen, I. Chumakov, and J. Weissenbach. A first-generation physical map of the human genome. *Nature*, 366(6456):698–701, Dec 1993.
- [CHL01] M. Crochemore, C. Hancart, and T. Lecroq. *Algorithmique du texte*. Vuibert, 2001.
- [DDR99] O. Delgrange, M. Dauchet, and E. Rivals. Location of Repetitive Regions in Sequences By Optimizing A Compression Method. In Russ Altman, editor, *Proc. of the 4th Pacific Symposium on Biocomputing*, volume 4, Hawaii, Jan 4-9 1999.
- [DF99] R. G. Downey and M. R. Fellows. *Parameterized Complexity*. Springer, 1999.
- [EG02] Olivier Elemento and Olivier Gascuel. An efficient and accurate distance based algorithm to reconstruct tandem duplication trees. In Thomas Lengauer and Hans-Peter Lenhof, editors, *Proceedings of the European Conference on Computational Biology 2002 (ECCB-02)*, volume 18, 2 of *Bioinformatics*, pages 92–99, London, October 6–9 2002. OXFORD University Press.
- [EG03] Olivier Elemento and Olivier Gascuel. An exact and polynomial distance-based algorithm to reconstruct single copy tandem duplication trees. In Ricardo Baeza-Yates and Maxime Crochemore, editors, *Combinatorial Pattern Matching, 14th Annual Symposium*, Lecture Notes in Computer Science. Springer, 2003.
- [EGL01] Olivier Elémento, Olivier Gascuel, and Marie-Paule Lefranc. Reconstruction de l’histoire de duplication de gènes répétés en tandem. In *Actes de JOBIM*, pages 9–11, Toulouse, France, 2001.
- [EGL02] Olivier Elémento, Olivier Gascuel, and Marie-Paule Lefranc. Reconstructing the duplication history of tandemly repeated genes. *Molecular Biology and Evolution*, 19(3):278–288, 2002.
- [Fit71] W. M. Fitch. Towards defining the course of evolution: minimum change for a specified tree topology. *Systematic Zoology*, 20:406–416, 1971.
- [Fit77] W. M. Fitch. Phylogenies constrained by the crossover process as illustrated by human hemoglobins and a thirteen-cycle, eleven-amino-acid repeat in human apolipoprotein A-I. *Genetics*, 86(3):623–44, Jul 1977.
- [Gas97] O. Gascuel. BIONJ: an improved version of the NJ algorithm based on a simple model of sequence data. *Molecular Biology and Evolution*, 14(7):685–95, Jul 1997.
- [GHJMM03] Olivier Gascuel, Michel D. Hendy, Alain Jean-Marie, and Robert McLachlan. The Combinatorics of Tandem Duplication Trees. *Systematic Biology*, 52(1):110–118, 2003.
- [GJW85] P. Gill, A. J. Jeffreys, and D. J. Werrett. Forensic application of DNA ‘fingerprints’. *Nature*, 318(6046):577–9, Dec 1985.
- [Gus97] Dan Gusfield. *Algorithms on Strings, Trees and Sequences*. Cambridge University Press, 1997.
- [HGH98] J. Hardy and K. Gwinn-Hardy. Genetic classification of primary neurodegenerative disease. *Science*, 282(5391):1075–9, Nov 1998.
- [HJ85] A. V. Hill and A. J. Jeffreys. Use of minisatellite DNA probes for determination of twin zygosity at birth. *Lancet*, 2(8469-70):1394–5, Dec 1985.
- [HP82] M. Hendy and D. Penny. Branch and bound algorithm to determine minimal evolutionary trees. *Mathematical Biosciences*, 59:277–290, 1982.
- [JAHS92] A. J. Jeffreys, M. J. Allen, E. Hagelberg, and A. Sonnberg. Identification of the skeletal remains of josef mengele by dna analysis. *Forensic Sci Int*, 56(1):65–76, 1992.
- [JBT98] M. A. Jobling, N. Bouzekri, and P. G. Taylor. Hypervariable digital DNA codes for human paternal lineages: MVR-PCR at the Y-specific minisatellite, MSY1 (DYF155S1). *Hum Mol Genet*, 7(4):643–53, 1998.

- [JKHM02] D. Jaitly, P.E. Kearney, G. Hui, and B. Ma. Methods for reconstructing the history of tandem repeats and their application to the human genome. *J. of Computer and System Sciences*, 2002.
- [JMT⁺91] A. J. Jeffreys, A. MacLeod, K. Tamaki, D. L. Neil, and D. G. Monckton. Minisatellite repeat coding as a digital approach to DNA typing. *Nature*, 354(6350):204–9, 1991.
- [JWT85a] A. J. Jeffreys, V. Wilson, and S. L. Thein. Hypervariable ‘minisatellite’ regions in human DNA. *Nature*, 314(6006):67–73, Mar 1985.
- [JWT85b] A. J. Jeffreys, V. Wilson, and S. L. Thein. Individual-specific ‘fingerprints’ of human DNA. *Nature*, 316(6023):76–9, Jul 1985.
- [KK00] R. Kolpakov and G. Kucherov. On maximal repetitions in words. *J. of Discrete Algorithms*, 1(1):159–186, 2000.
- [KK01] Roman Kolpakov and Gregory Kucherov. Finding approximate repetitions under Hamming distance. In *ESA: Annual European Symposium on Algorithms*, volume 2161 of *Lecture Notes in Computer Science*, pages 170–181. Springer, 2001.
- [Len02] Thomas Lengauer, editor. *Bioinformatics - From Genome to Drugs*, volume II: Applications of *Methods and Principles in Medicinal Chemistry*. Wiley-VCH Verlag, Weinheim, 2002.
- [Li97] Wen-Hsiung Li. *Molecular Evolution*. Sinauer Associates, Inc., Sunderland Massachusetts, U.S.A., 1997.
- [Lot99] M. Lothaire. *Algebraic Combinatorics on Words*. Cambridge University Press, 1999. URL: <http://www-igm.univ-mlv.fr/~berstel/Lothaire/index.html>.
- [Ohn70] S. Ohno. *Evolution by Gene Duplication*. Springer Verlag, New York, 1970.
- [PH98] Roderick D. M. Page and Edward C. Holmes. *Molecular Evolution: a Phylogenetic Approach*. Blackwell Science Ltd, Osney Mead, Oxford, 1998.
- [RDD⁺97] É. Rivals, O. Delgrange, J-P. Delahaye M. Dauchet, M-O. Delorme, A. Hénaut, and E. Ollivier. Detection of significant patterns by compression algorithms: the case of Approximate Tandem Repeats in DNA sequences. *Comp. Appl. in Biosciences*, 13(2):131–136, 1997.
- [RDDD96] É. Rivals, M. Dauchet, J-P. Delahaye, and O. Delgrange. Compression and genetic sequences analysis. *Biochimie*, 78(4):315–322, 1996.
- [SG02] J. Stoye and D. Gusfield. Simple and Flexible Detection of Contiguous Repeats Using a Suffix Tree. *Theoretical Computer Sciences*, 27(1-2):843–856, 2002.
- [SK99] David Sankoff and Joseph B. Kruskal, editors. *Time Warps, String Edits and Macromolecules: the Theory and Practice of Sequence Comparison*. CSLI Publications, second edition, 1999.
- [SM98] M. F. Sagot and E. W. Myers. Identifying satellites and periodic repetitions in biological sequences. *J. of Computational Biology*, 5(3):539–53, 1998.
- [TWY02] Mengxiang Tang, Michael Waterman, and Shibu Yooseph. Zinc Finger Gene Clusters and Tandem Gene Duplication. *J. of Computational Biology*, 9(2):429–446, 2002.
- [Wel96] R. D. Wells. Molecular basis of genetic instability of triplet repeats. *J Biol Chem*, 271(6):2875–8, Feb 1996.
- [WG97] Lusheng Wang and Dan Gusfield. Improved approximation algorithms for tree alignment. *Journal of Algorithms*, 25(2):255–273, November 1997.
- [WJG01] Lusheng Wang, Tao Jiang, and Dan Gusfield. A more efficient approximation scheme for tree alignment. *SIAM Journal on Computing*, 30(1):283–299, February 2001.
- [WJL96] Lusheng Wang, Tao Jiang, and E. L. Lawler. Approximation algorithms for tree alignment with a given phylogeny. *Algorithmica*, 16(3):302–315, September 1996.
- [ZMW02] Louxin Zhang, Bin Ma, and Lusheng Wang. Efficient methods for inferring tandem duplication history. *Lecture Notes in Computer Science*, 2452:97–111, 2002.

Acknowledgements

I would like to thank François Nicolas, Sèverine Bérard, Denis Bertrand, Sylvie Pinloche and Olivier Gascuel for their help and suggestions on the manuscript. This work is supported by: the French Inter-EPST program for Bioinformatics, the Génopole of Montpellier, the Specific Actions “Algorithms in Biology” and “Algorithms and Sequences” of the STIC section of CNRS.