

Projeto de pesquisa: Estudo de seqüências e aplicações à Biologia Computacional

Augusto Fernandes Vellozo

26 de outubro de 2006

Resumo

O interesse principal deste projeto é a elaboração de algoritmos eficientes que permitam tratar e extrair informações presentes em grandes volumes de dados como seqüências biológicas. Estruturas de dados apropriadas e algoritmos eficientes são particularmente importantes no estudo de seqüências biológicas devido à crescente quantidade de organismos sequenciados e ao grande comprimento destas seqüências.

Peça fundamental na comparação de seqüências biológicas é a obtenção de alinhamentos de seqüências, que no entanto costumam desprezar eventos mutacionais importantes como a duplicação e a inversão. Em nosso trabalho em Biologia Computacional temos obtido algoritmos ainda polinomiais exatos e que reintroduzem estes eventos em sua análise. Este projeto inclui avançar nos trabalhos já iniciados nestes assuntos, assim como trabalhar em novos problemas dentro de outros assuntos da Biologia Computacional, como pontos de quebra (*breakpoints*), motivos (*motifs*), genômica comparativa e técnicas para alinhamentos múltiplos.

O laboratório LBBE (*Laboratoire de Biométrie et Biologie Évolutive*) da Universidade Claude Bernardde - Lyon I ligado ao INRIA, onde será executado este plano de trabalho, assim como a supervisora Marie-France Sagot, têm sua excelência em pesquisa reconhecidas internacionalmente.

1 Introdução e justificativas

1.1 Comparação de seqüências

À medida que o número de novos genomas completos aumenta, a comparação entre seqüências longas de DNA de espécies próximas torna-se mais importante para nosso entendimento da estrutura da seqüência do DNA. Devido a isto, a análise genômica comparativa [30], apesar de ser um novo campo na bioinformática, está se desenvolvendo rapidamente. Em muitas espécies próximas, a ordem dos genes é preservada para intervalos curtos [35], i.e. sintenia. Nesses casos, os genes são mais conservados do que as regiões intergênicas. Portanto, a ordem da seqüência de genes é muito útil para detectar reordenamentos cromossômicos como inversões. Estes tipos de comparações ganham maior significância à medida que mais segmentos de genomas ortólogos, fortemente relacionados pela evolução, são sequenciados.

Desde a finalização do rascunho do genoma humano novos projetos de sequenciamento têm sido desenvolvidos para comparação com o genoma humano. Muitos programas computacionais têm sido usados para esse propósito como VISTA [33, 17], GLASS [7], Mummer [14], PipMaker [43], e também BLAST 2 Sequences [44].

Assim como muitos outros estudos em bioinformática, essas comparações dependem fortemente da obtenção de um alinhamento ótimo.

Na história da evolução alguns eventos introduzem mudanças na seqüência do DNA. Alguns eventos biológicos típicos são as *substituições*, *remoções* e *inserções* de nucleotídeos. Portanto, qualquer comparação de seqüências precisa levar em consideração a possibilidade da ocorrência desses eventos, se é esperado identificar uma alta similaridade entre duas seqüências. Procedimentos de alinhamento típicos tentam identificar que partes não mudam e onde se localizam esses eventos biológicos. Após, apresentam um alinhamento ótimo de acordo com algum critério de otimização e sistema de pontuação associado aos eventos.

Alinhamentos podem ser associados a um conjunto de operações de edição que transformam uma seqüência em outra. Normalmente as únicas operações de edição consideradas são a *substituição* de um símbolo em outro, a *inserção* de um símbolo e a *remoção* de um símbolo. Se os custos são associados a cada operação, existe um procedimento de programação dinâmica clássico em

$O(n^2)$ ¹ que computa o conjunto mínimo de operações de edição com o custo total mínimo e apresenta o alinhamento associado, que tem boa qualidade e alta semelhança para custos realistas.

1.2 Inversões

Considere três novas possibilidades de operações de edição:

- a *reversão-por-2*, que reverte a ordem de *dois símbolos consecutivos*;
- a operação de *reversão*, que reverte a ordem de *qualquer segmento* de símbolos ao invés de um segmento de comprimento 2;
- a operação de *inversão*, que substitui qualquer segmento pela sua sequência *reversa complementar*. A operação de inversão é a operação de interesse na obtenção de alinhamentos ótimos de seqüências biológicas.

Associados a quaisquer dessas três operações, nós podemos definir novos problemas de alinhamento. Por exemplo, dadas duas seqüências e custos fixos para cada tipo de operação de edição, o problema de *alinhamento com inversões* é um problema de otimização que indaga o custo total mínimo de um conjunto de operações de edição que transforma uma seqüência em outra. Além disso, pode também haver interesse na apresentação do alinhamento correspondente, ou seja das operações de edição. Da mesma forma, podemos definir os problemas de *alinhamento com reversões-por-2* e de *alinhamento com reversões*.

Em 1975, Wagner [48] estudou o alinhamento com reversões-por-2 e provou que ele admite uma solução polinomial se o custo de uma reversão-por-2 é nulo. Por outro lado, ele também provou que a obtenção de uma solução ótima é *NP-difícil*, se cada operação tem um custo positivo constante.

Em [11] podemos ver que o problema de decisão, associado ao problema de obter um alinhamento considerando inversões para um alfabeto de tamanho ilimitado, é NP-difícil.

Com o objetivo de tratar os alinhamentos com inversões, três estratégias principais têm sido consideradas:

- inversões sem sobreposições;
- ordenação de permutações sem sinal por reversões e;
- ordenação de permutações com sinal por reversões.

¹Consideraremos neste texto que n é o comprimento da maior seqüência analisada

Em 1992, Schöniger e Waterman [42] introduziram uma *hipótese simplificada*: *todas as regiões envolvendo inversões não se sobrepõem*. Isso levou ao problema do *alinhamento com inversões sem sobreposições* (primeira estratégia). Eles apresentaram uma solução em $O(n^6)$ para esse problema e também introduziram uma heurística que reduziu a complexidade do tempo de execução do problema. Essa heurística usa o algoritmo desenvolvido por Waterman e Eggert [49] que informa os K melhores alinhamentos locais não mutualmente intersectantes, com o objetivo de reduzir o tempo de execução para algo entre $O(n^2)$ e $O(n^4)$, dependendo dos dados.

A segunda estratégia se aplica bem a alinhamentos de *seqüências de genes* e tem sido bastante usada para genomas de mitocôndrias. Ela não se aplica a seqüências de nucleotídeos nem a seqüências de aminoácidos porque *repetições* de símbolos *não são* permitidas. Além disso, *nenhuma inserção* e *nenhuma remoção* são consideradas e a única operação permitida é a reversão em que a *reversão* é definida para transformar uma seqüência tal como 1, 2, 3, 4, 5 em 1, 4, 3, 2, 5. O problema, também chamado de *ordenação de permutações sem sinal por reversões*, produz a média do cálculo das distâncias de edição de duas permutações com a operação de reversão. Neste caso, os dados são duas permutações de $1, 2, 3, \dots, n$, onde n é o número de genes. Kececioglu e Sankoff [28] propuseram um algoritmo de 2-aproximação em 1995 e Christie [12] propôs um algoritmo de aproximação de razão $3/2$ em 1998. De fato, Caprara [10] provou em 1999 que esse problema na verdade é NP-difícil.

A terceira estratégia é o problema chamado *ordenação de permutações com sinal por reversões*. Este é o mesmo problema de ordenação de permutações sem sinal por reversões até o ponto em que os sinais também são atribuídos a um gene e uma reversão também troca seu sinal. Por exemplo, uma reversão poderia transformar 1, 2, 3, 4, 5 em 1, -4 , -3 , -2 , 5. Este sinal é normalmente associado à direção do gene (a qual filamento de DNA ele pertence). Hannenhalli e Pevzner [24] propuseram o primeiro algoritmo polinomial para o problema em 1995 e iniciaram uma seqüência de artigos baseados nessa estratégia. O algoritmo de Hannenhalli e Pevzner era $O(n^4)$ e foi melhorado para $O(n^2)$ por Kaplan, Shamir e Tarjan [26, 27] em 1997. Em 2001, Bader, Moret, e Yan [5] propuseram um algoritmo que calcula a distância de edição em $O(n)$ (a seqüência de reversões ainda requer $O(n^2)$). Estes estudos têm sido aplicados a estudos de reconstrução filogenética.

Em 2000, El-Mabrouk [19, 20] estudou a inclusão de duas operações: *inserções* e *remoções* de segmentos genéticos. Ela obteve resultados parciais e

propôs uma solução polinomial exata para um caso e uma heurística polinomial com um testador polinomial para otimalidade no outro caso. Símbolos repetidos ainda não são permitidos. Em 2002, El-Mabrouk [21] também obteve alguns resultados parciais ao considerar reversões e duplicações.

Em 2003, Lago [15] propôs dois algoritmos exatos para o problema do alinhamento com inversões sem sobreposições, ou seja, a primeira estratégia para tratar o problema do alinhamento com inversões. Um algoritmo é uma solução que executa em tempo $O(n^4)$ e que usa espaço $O(n^2)$. O outro algoritmo é uma implementação dinâmica esparsa que reduz o uso de recursos se $o(n^2)$ atribuições são dadas. Isto é freqüentemente esperado se a cardinalidade do alfabeto for grande, como por exemplo quando as letras são fragmentos de DNA de comprimento fixo.

Em 2005, Alves, Lago, e Vellozo [3] propusemos um algoritmo exato para este mesmo problema do alinhamento com inversões não sobrepostas que executa em tempo $O(n^3 \log n)$ e que usa espaço $O(n^2)$.

Em 2006, Vellozo, Alves, e Lago [46] propusemos um algoritmo exato para o mesmo problema que executa em tempo $O(n^3)$ e que usa espaço $O(n^2)$.

1.3 *Motifs* e repetições

Em genética, um motivo, ou *motif*, é um padrão de sequência de nucleotídeos ou aminoácidos que é amplamente encontrado e tem, ou espera-se que tenha, um significado biológico. Existem também os *motifs* estruturais, aplicados principalmente em proteínas e que estabelecem padrões tridimensionais.

Um exemplo de *motif* é o *N-glycosylation motif*, cujo padrão é descrito a seguir: *Asn*, seguido por qualquer aminoácido exceto *Pro*, seguido por *Ser* ou *Thr*, seguido por qualquer aminoácido exceto *Pro*, onde as abreviações de 3 letras identificam um tipo de aminoácido. Este padrão pode ser escrito como $N\{P\}[ST]\{P\}$, onde $N = Asn$, $P = Pro$, $S = Ser$, $T = Thr$; $\{X\}$ significa qualquer aminoácido exceto X ; e $[XY]$ significa ou X ou Y . A notação $[XY]$ não explicita a distribuição de probabilidade das ocorrências de X e de Y . Quando isto é desejado, estes padrões são definidos a partir de um modelo estatístico.

A procura de *motifs* nas seqüências normalmente está relacionada à busca de repetições nas seqüências.

Repetições são elementos preponderantes, especialmente em genomas de organismos eucariontes. Estima-se que mais de 80% dos genomas de planta são compostos por repetições. Existem diversos tipos de repetições em um

genoma, provavelmente nem todos já conhecidos. Entre os tipos mais conhecidos estão os satélites (micro ou mini conforme as características de comprimento e outras) que são repetições em *tandem*, isto é, que aparecem uma atrás da outra ao longo do genoma. Outro tipo de repetição muito conhecido são os ditos elementos transponíveis, ou transposons. Os transposons foram descobertos por Barbara McClintock [13] nos anos 50 estudando o milho. Os elementos transponíveis podem ser definidos como seqüências de DNA moderadamente repetitivas que podem mover-se de um local a outro no genoma e, desta maneira, ter um profundo impacto na estrutura, regulação e função dos genes, bem como na organização dos cromossomos na espécie. Enfim, motivos (*motifs*) em seqüências potencialmente relacionados com sítios de interação de complexos moleculares (proteínas e/ou RNAs) com o DNA são mais um exemplo de repetições, intra e inter espécie que desempenha um papel importante na regulação individual dos genes.

Ser capaz de identificar de maneira sistemática, e portanto exata, dado uma certa definição de repetição, é um problema importante em bioinformática que ainda não foi resolvido de maneira satisfatória ou realmente eficaz, especialmente no caso de genomas de organismos eucariontes. A dificuldade vem já da grande variedade de tipos de repetições. Algumas repetições, como as compridas (onde cada cópia pode atingir centenas de bases), são particularmente difíceis de identificar, por causa do comprimento, e porque há pouca conservação de uma cópia para a outra. Obviamente, o tamanho habitual das seqüências onde as repetições devem ser identificadas, aumenta ainda mais o grau de complexidade do problema.

Algumas doenças humanas são associadas às repetições, tais como: retardação mental *fragile-X* [47], doença de Huntington [1], distrofia miotônica [22] e ataxia de Friedreich [9]. *Tandem repeats* podem estar ligados a regras de regulação gênica [23, 32, 36], ligação DNA-proteína [39, 51] e evolução [25].

O número de cópias num *tandem repeat* pode ser variável entre indivíduos diferentes (polimórfico). Locais polimórficos são úteis em várias tarefas de laboratório [18, 50]. *Tandem repeats* tem sido utilizados para sustentar algumas hipóteses da evolução humana [4, 45] e da evolução de micro-satélites (*tandem repeats* cujo tamanho é de apenas algumas unidades de nucleotídeos) em primatas [34].

Em 1997, Benson [8] propôs um modelo para o alinhamento de seqüências que considera *tandem repeats* de mesmo tamanho. Ele propôs dois algoritmos exatos para obter um tal alinhamento ótimo. O primeiro algoritmo proposto

executa em tempo $O(n^5)$ e espaço $O(n^2)$. O segundo algoritmo proposto executa em tempo $O(n^4)$ e espaço $O(n^3)$.

Estamos elaborando um artigo com um algoritmo que obtém um alinhamento ótimo considerando o mesmo modelo proposto por Benson, porém com tempo de execução $O(n^3)$ e espaço $O(n^2)$.

1.4 Pontos de quebra

Transposons são uma espécie do que chamamos de "rearranjos genômicos". Outros rearranjos genômicos são a inversão de um trecho de sequência genômica, a duplicação, a deleção ou perda de uma parte de um genoma ou a operação inversa, isto é a inserção, e, no caso de genomas multicromossômicos, a translocação ou troca de material gênico entre dois cromossomos homólogos, a fusão de dois cromossomos em um ou a fissão de um cromossomo que dá nascimento a dois.

Rearranjos são eventos com consequências funcionais importantes para um organismo. São no entanto eventos cujos mecanismos conhecemos ainda pouco. No âmbito de entender melhor tais mecanismos, torna-se importante ser capaz de localizar precisamente o que se chama de pontos de quebra, isto é, os pontos exatos onde um genoma foi "quebrado" como consequência de um rearranjo. Não existe atualmente nenhum algoritmo que permita uma tal identificação com precisão suficiente para permitir o estudo dessas regiões e determinar características possíveis, no nível do genoma, dos pontos de quebra, e portanto dos mecanismos que levam a uma quebra.

Essa localização leva a problemas de alinhamento um pouco particulares, e que são especialmente difíceis, pois as sequências a serem alinhadas de uma maneira precisa são sequências frequentemente muito comprimidas, e com pouquíssima conservação. Trata-se, de fato, de sequências inter-gênicas onde se exerce, em geral, pouca pressão seletiva.

Essa representa, no entanto, somente uma primeira etapa para a compreensão dos possíveis mecanismos que podem levar a um rearranjo. Uma segunda etapa seria ser capaz de identificar características que seriam próprias a essas regiões que, em certos organismos, são "quebrados", e eventualmente até próprias ao tipo de rearranjo. Para tal, devemos novamente fazer comparação de sequências para tentar determinar tais características. A principal é que uma região que quebrou uma vez é possivelmente mais propensa a quebras múltiplas, e que as regiões ao redor de um ponto de quebra inicialmente identificado foram rearranjadas mais de uma vez. Compará-las sig-

nifica ser capaz de alinhar seqüências levando em conta diversos tipos mais de operações do que as usuais de substituição, deleção puntual, inserção e inversão (que eu fiz no meu doutorado).

A detecção exata dos *breakpoints* poderia ajudar a entender os mecanismos que estão por trás dos rearranjos e quais na verdade aconteceram. Isto poderia ajudar a melhorar os modelos de comparação de genomas e sua evolução. A hipótese de *Hotspots* (regiões que são mais suscetíveis a rearranjos [6, 29, 37, 38, 40]) é um assunto que poderia ser melhor estudado com a detecção exata dos *breakpoints*.

Em suma, a reconstrução do quebra-cabeças que é a evolução de um genoma leva a problemas algorítmicos de alinhamento e comparação que estão ainda largamente abertos.

2 Objetivos

Os objetivos deste projeto são relacionados a seguir:

- A elaboração de algoritmos eficientes que permitam tratar e extrair informações de seqüências biológicas;
- submeter ao menos dois artigos em publicações com reconhecimento internacional;
- possibilitar uma grande troca de conhecimentos científicos e acadêmicos entre os pesquisadores envolvidos;
- estabelecer contatos com outros pesquisadores da área de Biologia Computacional e até mesmo com biólogos que trabalhem com análise de seqüências;
- estreitar e fortalecer o relacionamento com a pesquisadora Marie-France Sagot e com outros pesquisadores do laboratório LBBE (*Laboratoire de Biométrie et Biologie Évolutive*) da Universidade Claude Bernard de Lyon I.

Já foram identificados alguns tópicos que são de interesse de ambas as partes e que iremos investigar durante o projeto. Estes tópicos são:

1. melhorias e mudanças nos algoritmos para alinhamento com inversões não sobrepostas e para alinhamento com repetições desenvolvidos durante o doutorado;

2. buscas de inversões e repetições ocorridas durante a evolução de espécies próximas com genomas sequenciados;
3. busca de pontos de quebra (*breakpoints*) ocorridos durante o processo evolutivo;
4. buscas de motivos (*motifs*);
5. técnicas para alinhamento múltiplo (várias seqüências).

3 Plano de trabalho e cronograma

Durante meu pos-doutorado, pretendo estudar esses problemas, me apoiando para isso nos trabalhos que fiz durante meu doutorado.

Pretendemos implementar os algoritmos que vierem a ser elaborados e realizar testes juntamente com biólogos com os quais meu orientador de pós-doutorado já trabalha. Esse contato com experimentalistas deve também ajudar a definir de maneira adequada os tipos de repetições que serão interessantes buscar em genomas completos de eucariontes, como o da mosca, do homem, do camundongo e do chimpanzé. Meu trabalho se fará também em colaboração com doutorandos do grupo com quem irei trabalhar.

Eventualmente, poderemos também investigar a possibilidade de inclusão de técnicas aceleradoras como o uso de árvores de sufixos ou de filtragem [16] nos problemas que vão ser estudados.

Estaremos receptivos a trabalhar também com problemas novos que vierem a surgir durante o projeto.

Para os resultados relevantes obtidos redigiremos os respectivos artigos e sempre que possível, enviaremos-os para publicações de reconhecimento internacional na área.

Pretendo executar os tópicos de interesse comum, já identificados e descritos anteriormente, na seguinte ordem e prazos:

1. melhorias e mudanças nos algoritmos do doutorado (2 meses);
2. buscas de inversões e repetições (3 meses);
3. busca de pontos de quebra (*breakpoints*) (2 meses);
4. buscas de motivos (*motifs*) (3 meses);
5. técnicas para alinhamento múltiplo (várias seqüências) (2 meses).

4 Síntese da bibliografia de referência

Várias referências que serão utilizadas neste projeto de pós-doutorado já foram citadas e de certa forma brevemente introduzidas na seção 1. Além destas existem outras que acreditamos serem úteis para o pós-doutorado. Algumas delas brevemente descrevemos a seguir:

- No trabalho de Jeanette Schmidt [41], descreve-se, dadas duas seqüências, como construir em tempo $O(n^2 \log n)$ uma estrutura de árvores binárias que permite obter, em tempo $O(\log n)$, o valor do alinhamento de qualquer fator de uma seqüência contra qualquer prefixo da outra.
- Já o trabalho de Landau e Ziv-Ukelson [31] descreve como podemos computar a pontuação do alinhamento de uma seqüência t contra um segmento comum y de diversas seqüências s_i , $1 \leq i \leq r$, uma única vez e aproveitar o valor desta computação no cálculo das pontuações dos alinhamentos de t contra cada s_i , $1 \leq i \leq r$.
- Por fim, o trabalho de Aggarwal *et al.* [2] descreve dois tipos de matrizes, monotônicas e totalmente monotônicas, e exhibe algoritmos para calcular os valores máximos de cada coluna para matrizes $n \times m$, monotônicas e totalmente monotônicas, em tempo, $O(m \log n)$ e $O(m)$, respectivamente.

5 Material e Métodos

As implementações dos algoritmos serão desenvolvidos em microcomputadores, e serão processados dados reais obtidos do sequenciamento de seqüências biológicas. Estas seqüências serão obtidas através de bancos de dados públicos ou de algum contacto com biólogos locais.

Quaisquer informações e artigos a serem utilizados na pesquisa serão obtidos preferencialmente através de meios eletrônicos em formato digital. Caso isto não seja possível então serão utilizados os periódicos e publicações arquivados na universidade.

6 Forma de análise dos resultados

Os resultados obtidos com os algoritmos aqui propostos serão comparados com aqueles obtidos de eventuais programas já disponíveis.

Também pretendemos obter retorno da relevância e acuracidade dos resultados obtidos a partir de diálogos com biólogos e pesquisadores locais.

Referências

- [1] A novel gene containing a trinucleotide repeat that is expanded and unstable on Huntington’s disease chromosomes. The Huntington’s Disease Collaborative Research Group. *Cell*, 72(6):971–983, Mar 1993.
- [2] A. Aggarwal, M. M. Klawe, S. Moran, P. Shor, and R. Wilber. Geometric applications of a matrix-searching algorithm. *Algorithmica*, 2(2):195–208, 1987.
- [3] C. E. R. Alves, A. P. do Lago, and A. F. Vellozo. Alignment with non-overlapping inversions in $O(n^3 \log n)$ -time. *Electronic Notes in Discrete Mathematics*, 19:365–371, Jun 2005.
- [4] Armour, Anttinen, May, Vega, Sajantila, Kidd, Kidd, Bertranpetit, Pääbo, and Jeffreys. Minisatellite diversity supports a recent African origin for modern humans. *Nat Genet*, 13(2):154–160, Jun 1996.
- [5] D. A. Bader, B. M. E. Moret, and M. Yan. A linear-time algorithm for computing inversion distance between signed permutations with an experimental study. In *Algorithms and data structures (Providence, RI, 2001)*, volume 2125 of *Lecture Notes in Comput. Sci.*, pages 365–376. Springer, Berlin, 2001.
- [6] J. A. Bailey, R. Baertsch, W. Kent, D. Haussler, and E. E. Eichler. Hotspots of mammalian chromosomal evolution. *Genome Biology*, 5:R23, 2004.
- [7] S. Batzoglou, L. Pachter, J. P. Mesirov, B. Berger, and E. S. Lander. Human and mouse gene structure: comparative analysis and application to exon prediction. *Genome Research*, 10(7):950–958, Jul 2000.
- [8] G. Benson. Sequence alignment with tandem duplication. In *RECOMB ’97: Proceedings of the first annual international conference on Computational molecular biology*, pages 27–36, New York, NY, USA, 1997. ACM Press.
- [9] V. Campuzano, L. Montermini, M. D. Molto, L. Pianese, M. Cossee, F. Cavalcanti, E. Monros, F. Rodius, F. Duclos, A. Monticelli, F. Zara, J. Canizares, H. Koutnikova, S. I. Bidichandani, C. Gellera, A. Brice, P. Trouillas, G. de Michele, A. Filla, R. de Frutos, F. Palau, P. I. Patel, S. di Donato, J.-L. Mandel, S. Coccozza, M. Koenig, and M. Pandolfo. Friedreich’s Ataxia: Autosomal Recessive Disease Caused by an Intronic GAA Triplet Repeat Expansion. *Science*, 271:1423–1427, Mar. 1996.

- [10] A. Caprara. Sorting permutations by reversals and Eulerian cycle decompositions. *SIAM J. Discrete Math.*, 12(1):91–110 (electronic), 1999.
- [11] X. Chen, J. Zheng, Z. Fu, P. Nan, Y. Zhong, S. Lonardi, and T. Jiang. Assignment of orthologous genes via genome rearrangement. *IEEE/ACM Trans. Comput. Biol. Bioinformatics*, 2(4):302–315, 2005.
- [12] D. A. Christie. A $3/2$ -approximation algorithm for sorting by reversals. In *Proceedings of the Ninth Annual ACM-SIAM Symposium on Discrete Algorithms (San Francisco, CA, 1998)*, pages 244–252, New York, 1998. ACM.
- [13] N. C. Comfort. From controlling elements to transposons: Barbara McClintock and the nobel prize. *Endeavour*, 25(3):127–130, September 2001.
- [14] A. L. Delcher, S. Kasif, R. D. Fleischmann, J. Peterson, O. White, and S. L. Salzberg. Alignment of whole genomes. *Nucleic Acids Research*, 27(11):2369–2376, Jun 1999.
- [15] A. P. do Lago, I. Muchnik, and C. Kulikowski. A sparse dynamic programming algorithm for alignment with non-overlapping inversions. *Theor. Inform. Appl.*, 39(1):175–189, 2005.
- [16] A. P. do Lago and I. Simon. *Tópicos em Algoritmos sobre Sequências*. IMPA, Rio de Janeiro, 2003. ISBN 85-244-0202-4.
- [17] I. Dubchak, M. Brudno, G. G. Loots, L. Pachter, C. Mayor, E. M. Rubin, and K. A. Frazer. Active conservation of noncoding sequences revealed by three-way species comparisons. *Genome Research*, 10(9):1304–1306, Sep 2000.
- [18] Edwards, Hammond, Jin, Caskey, and Chakraborty. Genetic variation at five trimeric and tetrameric tandem repeat loci in four human population groups. *Genomics*, 12(2):241–253, Feb 1992.
- [19] N. El-Mabrouk. Genome rearrangement by reversals and insertions/deletions of contiguous segments. In *Combinatorial pattern matching (Montreal, QC, 2000)*, volume 1848 of *Lecture Notes in Comput. Sci.*, pages 222–234. Springer, Berlin, 2000.
- [20] N. El-Mabrouk. Sorting signed permutations by reversals and insertions/deletions of contiguous segments. *J. Discrete Algorithms*, 1(1):105–121, 2000.
- [21] N. El-Mabrouk. Reconstructing an ancestral genome using minimum segments duplications and reversals. *J. Comput. System Sci.*, 65(3):442–464, 2002. Special issue on computational biology 2002.

- [22] Y. H. Fu, A. Pizzuti, R. G. Fenwick, Jr., J. King, S. Rajnarayan, P. W. Dunne, J. Dubel, G. A. Nasser, T. Ashizawa, P. de Jong, B. Wieringa, R. Korneluk, M. B. Perryman, H. F. Epstein, and C. T. Caskey. An Unstable Triplet Repeat in a Gene Related to Myotonic Muscular Dystrophy. *Science*, 255:1256–1258, Mar. 1992.
- [23] H. Hamada, M. Seidman, B. H. Howard, and C. M. Gorman. Enhanced gene expression by the poly(dT-dG).poly(dC-dA) sequence. *Mol. Cell. Biol.*, 4(12):2622–2630, 1984.
- [24] S. Hannenhalli and P. Pevzner. Transforming cabbage into turnip: polynomial algorithm for sorting signed permutations by reversals. In *STOC '95: Proceedings of the twenty-seventh annual ACM symposium on Theory of computing*, pages 178–189, New York, NY, USA, 1995. ACM Press.
- [25] Hellman, Steen, Sundvall, and Pettersson. A rapidly evolving region in the immunoglobulin heavy chain loci of rat and mouse: postulated role of (dC-dA)_n.(dG-dT)_n sequences. *Gene*, 68(1):93–9100, Aug 1988.
- [26] H. Kaplan, R. Shamir, and R. E. Tarjan. Faster and simpler algorithm for sorting signed permutations by reversals. In *Proceedings of the Eighth Annual ACM-SIAM Symposium on Discrete Algorithms (New Orleans, LA, 1997)*, pages 344–351, New York, 1997. ACM.
- [27] H. Kaplan, R. Shamir, and R. E. Tarjan. A faster and simpler algorithm for sorting signed permutations by reversals. *SIAM J. Comput.*, 29(3):880–892 (electronic), 2000.
- [28] J. Kececioğlu and D. Sankoff. Exact and approximation algorithms for sorting by reversals, with application to genome rearrangement. *Algorithmica*, 13(1-2):180–210, 1995.
- [29] W. J. Kent, R. Baertsch, A. Hinrichs, W. Miller, and D. Haussler. Evolution’s cauldron: Duplication, deletion, and rearrangement in the mouse and human genomes. *Proceedings of the National Academy of Science*, 100:11484–11489, Sept. 2003.
- [30] E. V. Koonin. The emerging paradigm and open problems in comparative genomics. *Bioinformatics*, 15(4):265–266, Apr 1999. Editorial.
- [31] G. M. Landau and M. Ziv-Ukelson. On the common substring alignment problem. *J. Algorithms*, 41(2):338–359, 2001.

- [32] Q. Lu, L. L. Wallrath, H. Granok, and S. C. Elgin. (CT)_n (GA)_n repeats and heat shock elements have distinct roles in chromatin structure and transcriptional activation of the *Drosophila* hsp26 gene. *Mol. Cell. Biol.*, 13(5):2802–2814, 1993.
- [33] C. Mayor, M. Brudno, J. R. Schwartz, A. Poliakov, E. M. Rubin, K. A. Frazer, L. S. Pachter, and I. Dubchak. VISTA : visualizing global DNA sequence alignments of arbitrary length. *Bioinformatics*, 16(11):1046–1047, Nov 2000.
- [34] W. Messier, S. H. Li, and C. B. Stewart. The birth of microsatellites. *Nature*, 381:483, Jun 1996.
- [35] G. Moore, K. M. Devos, Z. Wang, and M. D. Gale. Cereal genome evolution. Grasses, line up and form a circle. *Current Biology*, 5(7):737–739, Jul 1995. Review.
- [36] Pardue, Lowenhaupt, Rich, and Nordheim. (dC-dA)_n(dG-dT)_n sequences have evolutionarily conserved chromosomal locations in *Drosophila* with implications for roles in chromosome structure and function. *EMBO J*, 6(6):1781–1789, Jun 1987.
- [37] Q. Peng, P. A. Pevzner, and G. Tesler. The Fragile Breakage versus Random Breakage Models of Chromosome Evolution. *PLoS Computational Biology*, 2:e14, Feb. 2006.
- [38] P. Pevzner and G. Tesler. Human and mouse genomic sequences reveal extensive breakpoint reuse in mammalian evolution. *Proceedings of the National Academy of Science*, 100:7672–7677, June 2003.
- [39] R. Richards, K. Holman, S. Yu, and G. Sutherland. Fragile X syndrome unstable element, p(CCG)_n, and other simple tandem repeat sequences are binding sites for specific nuclear proteins. *Hum. Mol. Genet.*, 2(9):1429–1435, 1993.
- [40] D. Sankoff. Chromosomal breakpoint reuse in genome sequence rearrangement. *Journal of Computational Biology*, 12(6):812–821, 2005.
- [41] J. P. Schmidt. All highest scoring paths in weighted grid graphs and their application to finding all approximate repeats in strings. *SIAM J. Comput.*, 27(4):972–992 (electronic), 1998.
- [42] M. Schöniger and M. S. Waterman. A local algorithm for DNA sequence alignment with inversions. *Bulletin of Mathematical Biology*, 54(4):521–536, Jul 1992.

- [43] S. Schwartz, Z. Zhang, K. A. Frazer, A. Smit, C. Riemer, J. Bouck, R. Gibbs, R. Hardison, and W. Miller. PipMaker—a web server for aligning two genomic DNA sequences. *Genome Research*, 10(4):577–586, Apr 2000.
- [44] T. A. Tatusova and T. L. Madden. BLAST 2 Sequences, a new tool for comparing protein and nucleotide sequences. *FEMS Microbiology Letters*, 174(2):247–250, May 1999.
- [45] S. A. Tishkoff, E. Dietzsch, W. Speed, A. J. Pakstis, J. R. Kidd, K. Cheung, B. Bonne-Tamir, A. S. Santachiara-Benerecetti, P. Moral, M. Krings, S. Pabo, E. Watson, N. Risch, T. Jenkins, and K. K. Kidd. Global Patterns of Linkage Disequilibrium at the CD4 Locus and Modern Human Origins. *Science*, 271:1380–1387, Mar. 1996.
- [46] A. F. Vellozo, C. E. R. Alves, and A. P. do Lago. Alignment with non-overlapping inversions in $o(n^3)$ -time. In *6th Workshop on Algorithms in Bioinformatics*. Springer, 2006. Lecture Notes in Bioinformatics 4175.
- [47] Verkerk, Pieretti, Sutcliffe, Fu, Kuhl, Pizzuti, Reiner, Richards, Victoria, and Zhang. Identification of a gene (FMR-1) containing a CGG repeat coincident with a breakpoint cluster region exhibiting length variation in fragile X syndrome. *Cell*, 65(5):905–914, May 1991.
- [48] R. Wagner. On the complexity of the extended string-to-string correction problem. In *Seventh ACM Symposium on the Theory of Computation*. Association for Computing Machinery, 1975.
- [49] Waterman and Eggert. A new algorithm for best subsequence alignments with application to tRNA-rRNA comparisons. *Journal of Molecular Biology*, 197(4):723–728, Oct 1987.
- [50] Weber and May. Abundant class of human DNA polymorphisms which can be typed using the polymerase chain reaction. *Am J Hum Genet*, 44(3):388–396, Mar 1989.
- [51] H. Yee, A. Wong, J. van de Sande, and J. Rattner. Identification of novel single-stranded d(TC)_n binding proteins in several mammalian species. *Nucl. Acids Res.*, 19(4):949–953, 1991.