



Universidade Federal de Pernambuco
Centro de Informática

Graduação em Engenharia da Computação

**Construção online de wavelet tree no
formato da árvore de Huffman**

Ícaro Julião Monteiro Guerra

Trabalho de Graduação

Recife
25 de setembro de 2023

Universidade Federal de Pernambuco
Centro de Informática

Ícaro Julião Monteiro Guerra

Construção online de wavelet tree no formato da árvore de Huffman

Trabalho apresentado ao Programa de Graduação em Engenharia da Computação do Centro de Informática da Universidade Federal de Pernambuco como requisito parcial para obtenção do grau de Bacharel em Engenharia da Computação.

Orientador: *Paulo Gustavo Soares da Fonseca*

Recife
25 de setembro de 2023

Ficha de identificação da obra elaborada pelo autor,
através do programa de geração automática do SIB/UFPE

Monteiro Guerra, Ícaro Julião.

Construção online de wavelet tree no formato da árvore de Huffman / Ícaro
Julião Monteiro Guerra. - Recife, 2023.

30 p. : il., tab.

Orientador(a): Paulo Gustavo Soares da Fonseca

Trabalho de Conclusão de Curso (Graduação) - Universidade Federal de
Pernambuco, Centro de Informática, Engenharia da Computação - Bacharelado,
2023.

1. Wavelet tree. 2. Código de Huffman. 3. Algoritmo dinâmico. I. Soares da
Fonseca, Paulo Gustavo. (Orientação). II. Título.

000 CDD (22.ed.)

Agradecimentos

Agradeço primeiramente aos meus pais, por se dedicarem durante todos esses anos para que pudesse usufruir de uma educação de qualidade, que me proporcionou ingressar na Universidade e concluir a graduação com êxito.

Agradeço à Aline, minha parceira de vida, que me acompanhou nos altos e baixos desse trabalho, sempre me motivando a acreditar no meu sucesso.

Agradeço ao meu orientador, Paulo Fonseca, por se mostra sempre solícito para dar as coordenadas e ajudar a concluir este último desafio.

E por último, quero agradecer a mim mesmo, por nunca ter desistido apesar das adversidades, pois todo o esforço está sendo colhido agora.

Resumo

A *wavelet tree* é uma estrutura de dados bastante utilizada para representar longas cadeias de caracteres de forma otimizada e permitir consultas de *rank* e *select*. Alguns estudos já exploraram como modificar o formato dessa árvore para melhorar a complexidade de consultas e de espaço de armazenamento e mostram que o formato da árvore de Huffman é ótimo no quesito memória. Apesar de ser uma estrutura conhecida, poucos algoritmos de construções *online* foram propostos e este trabalho tem como objetivo propor um algoritmo *online* para a construção da *wavelet tree* no formato da árvore de Huffman, que tem como principal vantagem o fato de que o texto de referência não precisa ser armazenado em memória principal, o que pode ser um requisito desejável, ou mesmo necessário, em aplicações de *streams* de dados. O método proposto consiste na adaptação de um algoritmo para construção online do Código de Huffman, representado por uma árvore estritamente binária com topologia variável. A mudança na forma da árvore implica em modificações nos bitvectors representados pelos nós, com impacto significativo no desempenho do algoritmo. O método foi implementado e testado em textos de diferentes características, e os resultados sugerem que pode representar uma alternativa viável, especialmente em memória, porém previsivelmente menos eficiente em tempo com respeito à construção offline tradicional.

Palavras-chave: Wavelet tree, árvore de Huffman, dinâmico, online

Sumário

1	Introdução	1
2	Métodos	3
2.1	Codificação de Huffman	3
2.1.1	Construção da árvore de Huffman dinâmica	6
2.2	Wavelet Tree	9
2.2.1	Construção dinâmica da <i>wavelet tree</i> no formato da árvore de Huffman	12
2.2.1.1	Remoção de <i>bits</i>	16
2.2.1.2	Adição de <i>bits</i>	17
2.3	Implementação	19
3	Resultados	21
3.1	Dados de teste	21
3.2	Avaliação	21
4	Conclusões	27

Lista de Figuras

2.1	Exemplo da árvore de Huffman estática construída sobre <i>abracadabra</i> .	7
2.2	Exemplo da transformação da árvore de Huffman no processamento da palavra <i>abra</i> para <i>abrac</i> . Os nós trocados estão marcados de cinza e W é a lista de pesos ordenados da propriedade dos irmãos.	10
2.3	A <i>wavelet tree</i> balanceada do texto <i>abracadabra</i> sobre o alfabeto a, b, c, d, r . Apenas os <i>bitvectors</i> são armazenados. O texto e alfabeto são para fins ilustrativos.	12
2.4	A <i>wavelet tree</i> no formato Huffman do texto <i>abracadabra</i> sobre o alfabeto a, b, c, d, r . Apenas o código de Huffman de cada caractere e os <i>bitvectors</i> são armazenados. O texto e alfabeto são para fins ilustrativos.	13
2.5	Exemplo quando o vértice v não será modificado na troca dos vértices x_i e x_j , onde $LCA(x_i, x_j) = x_k$.	15
2.6	Exemplo quando o vértice v será modificado na troca dos vértices x_i e x_j , onde $LCA(x_i, x_j) = x_k$ e $x_k = v$.	15
2.7	Exemplo quando o vértice v será modificado na troca dos vértices x_i e x_j , onde $LCA(x_i, x_j) = x_k$.	16
2.8	Exemplo de valores de P na remoção dos <i>bits</i> de x_i na troca com x_j . Os <i>bits</i> em negrito serão removidos.	18
2.9	Exemplo de valores de L na adição dos <i>bits</i> de x_j na troca com x_i . Os <i>bits</i> em negrito serão adicionados.	19
3.1	Tempo de processamento para cada tamanho de texto de inglês, DNA e proteína.	22
3.2	Número de troca de <i>bits</i> para cada tamanho de texto de inglês, DNA e proteína.	24

Lista de Tabelas

3.1	Métricas para textos em inglês de diferentes tamanhos	22
3.2	Métricas para sequências de DNA de diferentes tamanhos	23
3.3	Métricas para sequências de proteínas de diferentes tamanhos	24

Introdução

1

2 A representação de textos de forma eficiente é um dos problemas mais clássicos na computação.
3 Um dos algoritmos mais famosos e utilizados para representar dados de maneira eficiente é a
4 codificação de Huffman, uma técnica proposta por Huffman [9] para determinar o código de
5 prefixo de custo mínimo de um texto. Na construção desse código, é feita a construção de uma
6 estrutura auxiliar chamada árvore de Huffman que é utilizada para descobrir o código de cada
7 caractere do texto.

8 Nos anos 70, Faller [4] e Gallager [6] desenvolveram independentemente um algoritmo de
9 compressão dinâmico baseado no método de Huffman. No livro de Crochemore e Rytter [1],
10 temos a descrição do algoritmo, que gera o código de um texto sem conhecimento prévio do
11 seu alfabeto e que pode ter um uso de memória menor que o algoritmo proposto inicialmente
12 se levarmos em conta que não é necessário armazenar o texto para processamento.

13 A *wavelet tree* é uma estrutura de dados sucinta proposta por Grossi, Grupta e Vitter [7]
14 com uso em diversas áreas [12, 5]. Essa estrutura é utilizada para representar uma grande
15 quantidade de caracteres com em espaço comprimido e suporta responder diferentes tipos de
16 consulta em seu conteúdo de forma eficiente, como por exemplo saber quantos caracteres de
17 um certo tipo existe entre duas posições do texto. Por ser inicialmente proposta como uma
18 árvore binária balanceada, Grossi et al. [8] nota que as consultas funcionariam independente
19 do formato da árvore e propõe mudanças no formato para otimizar as consultas. Já Mäkinen
20 e Navarro [10], propõem a modificação do formato da árvore para o mesmo formato da árvore
21 de Huffman para conseguir uma compressão de espaço ainda melhor.

22 O trabalho de Fonseca e Silva [2] propõe um algoritmo para a construção *online* da *wavelet*

1 *tree* no formato de uma árvore binária balanceada que possui performance similar em questão
2 de tempo de processamento comparado aos algoritmos mais utilizados, mas uma melhoria no
3 espaço utilizado por não precisar armazenar o texto como um todo e utilizar pouca memória
4 extra além do necessário para a estrutura.

5 Este trabalho tem como objetivo propor um algoritmo de construção dinâmico da *wavelet*
6 *tree* com o formato da árvore de Huffman, não necessitando armazenar o texto em memória
7 nem saber o alfabeto com antecedência.

Métodos

1

2 Neste capítulo, detalhamos o algoritmo para construção de uma *Wavelet tree* dinâmica no for-
3 mato da árvore da codificação de Huffman. Primeiramente é necessário entender como funci-
4 ona a codificação de Huffman e sua construção dinâmica, para então modificarmos essa cons-
5 trução para incluir as informações necessárias para formação da *wavelet tree*.

2.1 Codificação de Huffman

6

7 A codificação de Huffman é uma técnica de compactação que tem como entrada um texto T ,
8 um alfabeto A de tamanho m , uma lista de pesos positivos $F = (f_i > 0 \mid 0 \leq i < m)$, em que f_i
9 é igual à frequência do caractere a_i em T , e tem como saída um código binário de prefixo de
10 custo mínimo, chamado de código de Huffman. Para entender o que isso significa, definimos:

- 11 1. Um código é binário quando o alfabeto de saída é igual à $\{0, 1\}$. Isso significa que existe
12 uma lista de valores $R = (r_i \in \{0, 1\}^* \mid 0 \leq i < m)$ onde r_i é a representação, ou também
13 chamado de código, do símbolo a_i no alfabeto de saída.
- 14 2. É considerado um código de prefixo quando o código de cada símbolo de A em R não é
15 um prefixo de nenhum outro. Por exemplo, com $m = 5$, uma possível representação do
16 alfabeto para ser um código de prefixo seria $R = \{00, 01, 10, 110, 111\}$, enquanto que a
17 representação $R = \{00, 01, 10, 11, 001\}$ é inválida porque 00 é prefixo de 001. Códigos
18 de prefixos também têm a característica de que são unicamente decodificados.

3. Um código é de custo mínimo quando utiliza a representação R que minimiza

$$C(R, F) = \sum_{i=0}^{m-1} \text{len}(r_i) \cdot f_i, \quad (2.1)$$

dentre todas as representações possíveis, onde $\text{len}(x)$ é igual ao tamanho ou quantidade de *bits* da representação x .

Para gerar a codificação de Huffman do texto T , é necessário construir a árvore de Huffman. A árvore de Huffman é uma árvore estritamente binária em que cada símbolo a_i corresponde a um nó folha de peso f_i , todo nó interno da árvore tem peso igual a soma do peso de seus dois filhos, cada aresta possui um símbolo de $\{0, 1\}$ e o código r_i pode ser obtido pelo caminho da raiz até o nó folha de a_i , concatenando os símbolos presentes nas arestas do caminho. É importante notar também que duas arestas que conectam ao mesmo nó pai não podem ter o mesmo símbolo. Sem perda de generalidade, assumimos que a aresta da subárvore esquerda terá o símbolo 0, e a aresta da subárvore direita terá símbolo 1.

Essa árvore pode ser construída seguindo um algoritmo guloso simples. Inicialmente, utilizamos cada símbolo de A e seu respectivo peso em F para construir um conjunto de nós folhas. Então, realizamos um processo guloso, onde os dois nós de menor peso são retirados do conjunto e um novo nó interno é criado, que possui peso igual a soma dos pesos dos dois nós retirados. Esse novo nó interno tem a aresta da subárvore esquerda conectada ao nó retirado de menor peso, e a aresta da subárvore da direita conectada ao outro nó retirado e esse novo nó é então adicionado de volta ao conjunto. Após $n - 1$ passos apenas um nó sobrar no conjunto, a raiz de nossa árvore de Huffman, com peso igual a $\sum_{i=0}^{n-1} f_i$. Um exemplo para a árvore de Huffman da palavra *abracadabra* pode ser visto na Figura 2.1.

Para codificar um texto T de tamanho n utilizando a árvore de Huffman, criamos a lista de códigos R a partir das folhas das árvores, e depois simplesmente concatenamos os códigos dos caracteres de T . Para $T = \textit{abracadabra}$, com $A = \{a, b, c, d, r\}$, usando a árvore da Figura 2.1 temos que $R = \{0, 10, 1100, 1101, 111\}$, e como resultado o código de Huffman $H(T) = 0 \cdot 10 \cdot$

1 111·0·1100·0·1101·0·10·111·0, isto é, $H(T) = 01011101100011010101110$.

2 A decodificação de um código de Huffman é simples: criamos um ponteiro, inicialmente
3 na raiz da árvore de Huffman, e processamos o código símbolo por símbolo. Para cada novo
4 símbolo do código, o ponteiro desce na árvore utilizando a aresta que possui mesmo valor.
5 Quando esse ponteiro alcança um nó folha, com um caractere e um peso associado, temos a
6 certeza que esse caractere é o caractere da posição atual e o ponteiro volta para a raiz. Esse
7 processo é realizado até o fim do código, e é válido porque código de Huffman é um código
8 de prefixo. Decodificando o código do exemplo anterior, $H = 01011101100011010101110$ e
9 utilizando a Figura 2.1 de referência, iniciamos o ponteiro na raiz e realizamos os seguintes
10 passos:

- 11 1. $h_0 = 0$, então o ponteiro desce para a esquerda. Como chegamos em um nó folha com o
12 ponteiro, temos que $T_0 = a$. Ponteiro volta para raiz.
- 13 2. $h_1 = 1$, ponteiro desce para a direita.
- 14 3. $h_2 = 0$, ponteiro desce para esquerda e $T_1 = b$. Ponteiro volta para raiz.
- 15 4. $h_3 = 1$, ponteiro desce pra direita.
- 16 5. $h_4 = 1$, ponteiro desce pra direita.
- 17 6. $h_5 = 1$, ponteiro desce pra direita e $T_2 = r$. Ponteiro volta para raiz.
- 18 7. $h_6 = 0$, ponteiro desce pra esquerda e $T_3 = a$. Ponteiro volta para raiz.
- 19 8. $h_7 = 1$, ponteiro desce pra direita.
- 20 9. $h_8 = 1$, ponteiro desce pra direita.
- 21 10. $h_9 = 0$, ponteiro desce pra esquerda.
- 22 11. $h_{10} = 0$, ponteiro desce pra esquerda e $T_4 = c$. Ponteiro volta para raiz.

- 1 12. $h_{11} = 0$, ponteiro desce pra esquerda e $T_5 = a$. Ponteiro volta para raiz.
 - 2 13. $h_{12} = 1$, ponteiro desce pra direita.
 - 3 14. $h_{13} = 1$, ponteiro desce pra direita.
 - 4 15. $h_{14} = 0$, ponteiro desce pra esquerda.
 - 5 16. $h_{15} = 1$, ponteiro desce pra direita e $T_6 = d$. Ponteiro volta para raiz.
 - 6 17. $h_{16} = 0$, ponteiro desce pra esquerda e $T_7 = a$. Ponteiro volta para raiz.
 - 7 18. $h_{17} = 1$, ponteiro desce pra direita.
 - 8 19. $h_{18} = 0$, ponteiro desce pra esquerda e $T_8 = b$. Ponteiro volta para raiz.
 - 9 20. $h_{19} = 1$, ponteiro desce pra direita.
 - 10 21. $h_{20} = 1$, ponteiro desce pra direita.
 - 11 22. $h_{21} = 1$, ponteiro desce pra direita e $T_9 = r$. Ponteiro volta para raiz.
 - 12 23. $h_{22} = 0$, ponteiro desce pra esquerda e $T_{10} = a$. Ponteiro volta para raiz.
 - 13 24. Fim do código H .
- 14 No final, temos que $T = abracadabra$, conforme esperado.

15 2.1.1 Construção da árvore de Huffman dinâmica

16 Descrevemos como funciona a construção da árvore de Huffman de maneira estática, ou seja:
17 tendo acesso ao texto T inteiro e calculando a frequência de cada caractere antes do início da
18 construção. Porém, aplicações que utilizam *stream* de dados podem não possuir acesso ao texto
19 por completo, e por isso se vê necessário o uso de um algoritmo dinâmico para construção da
20 árvore de Huffman, descrito em [1].

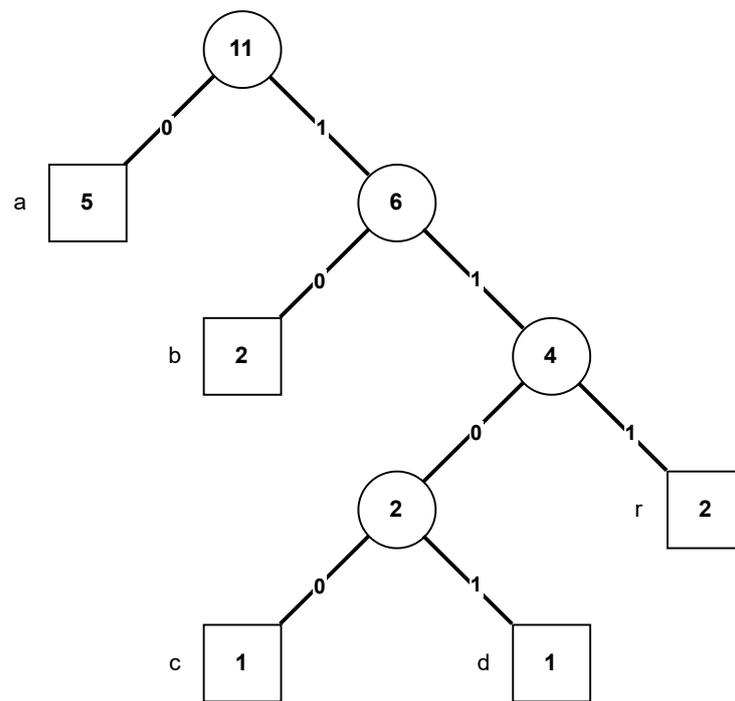


Figura 2.1 Exemplo da árvore de Huffman estática construída sobre *abracadabra*.

1 **Teorema 2.1** (*Propriedade dos irmãos*). *Seja V uma árvore binária completa, F uma lista de*
 2 *pesos positivos para cada nó da árvore, e p a quantidade de folhas da árvore, onde o peso de*
 3 *todo nó interno é igual a soma dos pesos de seus filhos. Então V é uma árvore de Huffman se*
 4 *somente se seus nós podem ser organizados numa sequência $S = (x_1, x_2, \dots, x_{2p-1})$ de nós da*
 5 *árvore tal que:*

- 6 1. *A sequência de pesos $(f_{x_1}, f_{x_2}, \dots, f_{x_{2p-1}})$ está em ordem crescente e*
- 7 2. *Para qualquer i ($1 \leq i < p$), os nós consecutivos x_{2i-1} e x_{2i} são irmãos (possuem o mesmo*
 8 *pai)*

9 *A prova pode ser vista em [1]*

10 A ideia do algoritmo dinâmico de Huffman é atualizar os nós da árvore de Huffman para
 11 cada mudança de frequência de um caractere de forma que a propriedade descrita no Teo-
 12 rema 2.1 se mantenha. Pode ser necessária tanto a criação de novos vértices, no caso de um
 13 caractere que apareceu pela primeira vez, como também a troca de posição entre nós. Pelas mu-
 14 danças de nó na árvore, o código de um caractere pode mudar ao longo do processamento de
 15 um mesmo texto e, portanto, se vê necessário a reconstrução da árvore com o mesmo algoritmo
 16 no momento de decodificação.

17 A árvore de Huffman dinâmica possui um nó a mais que a árvore de Huffman estática: um
 18 nó de peso 0 que será representado por um caractere que não aparecerá em T - no nosso caso
 19 utilizaremos #. Esse será o primeiro nó da árvore e único nó de peso 0. Pelo seu peso ser 0, ele
 20 sempre será o primeiro da sequência S do Teorema 2.1 e todos os novos nós folha criados serão
 21 vizinhos à ele inicialmente. Outra diferença do algoritmo estático é que precisamos definir uma
 22 função $g(c)$, que retorna um código binário único para todos os caracteres em A . Esse código
 23 será usado apenas na primeira aparição de c no texto.

24 Seja $T[l : r] = t_l \cdots t_r$ e suponha que processamos o prefixo $T[0 : i]$. Temos V^i , a árvore de
 25 Huffman para esse prefixo mais um nó folha de peso 0 para o caractere # que representa todos

1 os caracteres ainda não vistos e um código H^i . Ao ler o próximo caractere $c = T[i + 1]$, temos
 2 duas possibilidades:

- 3 1. c apareceu anteriormente em $T[0 : i]$. O código de c é encontrado igual a versão estática
 4 da árvore de Huffman: o valor das arestas no caminho da raiz até a folha.
- 5 2. c nunca apareceu antes em $T[0 : i]$, então não possui uma folha em V^i . Como essa é a
 6 primeira vez que c aparece no texto e não possui um código próprio ainda, utilizaremos
 7 o código de # concatenado com $g(c)$.

8 Após conseguir o código atual de c , atualizaremos V^i para V^{i+1} e teremos um H^{i+1} para $T[:$
 9 $i + 1]$. Com isso, podemos descrever como atualizar V^i para V^{i+1} . Caso c não tenha uma folha
 10 em V , iremos criar uma nova folha x_{2p} de peso 1 para representar c e um nó interno x_{2p+1} que
 11 será pai de x_{2p} e x_0 . Essa operação é igual a modificar a sequência S de $(x_0, x_2, \dots, x_{2p-1})$ para
 12 $(x_0, x_{2p}, x_{2p+1}, x_2, \dots, x_{2p-1})$ e isso mantém a propriedade dos irmãos válida para a sequência,
 13 já que $f_{x_0} = 0$ e $f_{x_{2p}} = f_{x_{2p+1}} = 1$. Por outro lado, se c já faz parte da árvore, começamos
 14 incrementando em 1 o peso f_i correspondente ao nó folha que representa c . Em ambos os
 15 casos, se a propriedade 1 do Teorema 2.1 for quebrada com esse incremento, procuraremos o
 16 maior j tal que $f_{s_i} > f_{s_j}$, $i < j$ e x_j não é pai de x_i e faremos a troca de posição desses nós,
 17 caso exista tal j . Esse processo é repetido sucessivamente para o pai de x_i , que tem seu peso
 18 incrementado em 1, até a raiz e pode ser implementado em $O(p)$, que também significa dizer
 19 que será $O(A)$. Ao final desse processo, teremos V^{i+1} . Um exemplo da troca de nós pode ser
 20 visto na Figura 2.2.

21 2.2 Wavelet Tree

22 Seja T um texto de tamanho n sobre um alfabeto A de tamanho m . Para todo subconjunto
 23 $A' \subseteq A$, seja $proj(T, A')$ a sequência não contínua, possivelmente vazia, de todos caracteres

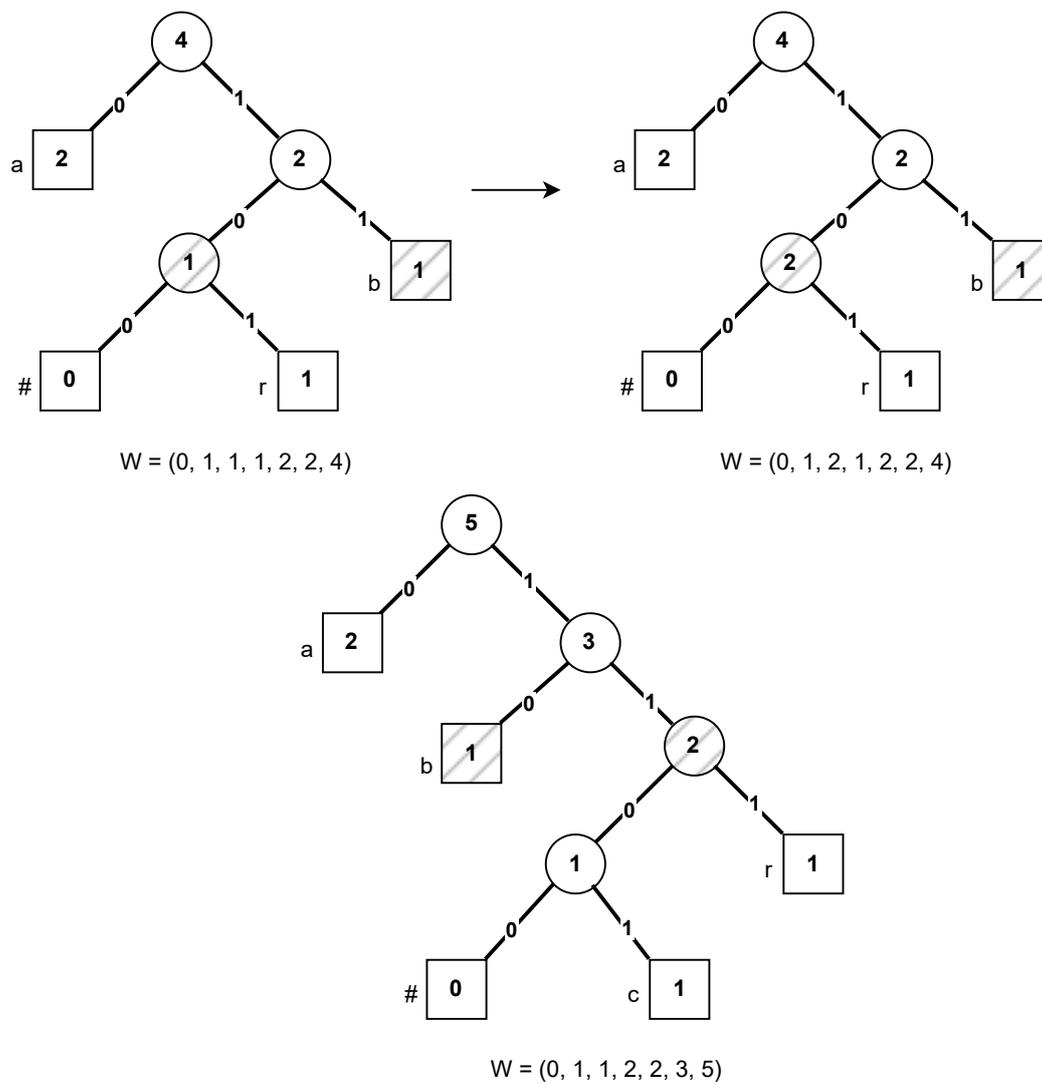


Figura 2.2 Exemplo da transformação da árvore de Huffman no processamento da palavra *abra* para *abrac*. Os nós trocados estão marcados de cinza e W é a lista de pesos ordenados da propriedade dos irmãos.

1 $t_i \in A'$. Por exemplo, seja $T = abracadabra$ e $A' = \{a, c, d\}$, $proj(T, A') = aacadaa$.

2 **Definição 2.1.** A *wavelet tree* (WT) para o texto T sobre o alfabeto A pode ser definida recur-
3 sivamente como

- 4 • Caso $|A| = 1$, então a árvore é apenas um nó folha de rótulo a_0 .
- 5 • Se não, a árvore possui um nó raiz interno, que representa A . Esse nó possui dois filhos:
6 o filho a esquerda é definido como a WT para o texto $proj(T, A_0)$ sobre o alfabeto A_0 e o
7 filho a direita é definido como a WT para o texto $proj(T, A_1)$ sobre o alfabeto A_1 , onde
8 $A_0 \cup A_1 = A$. Também armazenamos no nó um *bitvector* B de tamanho $|T|$, onde $b_i = 0$
9 se $t_i \in A_0$ ou $b_i = 1$ se $t_i \in A_1$.

10 A WT é uma estrutura que provê consultas de *access*, *rank* e *select*. Seja $W(T, A)$ a WT
11 sobre T em A e $T[l, r] = t_l \cdots t_r$, definimos que $W.access(i)$ é igual t_i , $W.rank(c, i)$ é igual ao
12 número de ocorrências de c em $T[0 : i]$ e $W.select(c, j)$ é igual à posição em T que o caractere
13 c aparece pela j -ésima vez. Então, se $i = W.select(c, j)$, temos $W.rank(c, i) = j$.

14 Na construção da WT, podemos definir diferentes formas de particionar A . Uma das formas
15 mais utilizadas é a de dividir A na metade a cada iteração. Essa forma de particionar A gera uma
16 WT balanceada, que, como demonstrado em [11], pode utilizar de $n \lg(m) + o(n \lg(m))$ bits de
17 memória e complexidade de consulta $O(\lg(m))$. Outra forma de particionar A , explorado em
18 [10], é a de dividir A de forma que a árvore resultante tenha o mesmo formato da árvore de
19 Huffman. É importante notar que, para essa divisão do alfabeto, é também necessário armaze-
20 nar em cada nó interno uma lista D , onde, para todo $c \in A$, $d_c = 0$ se $c \in A_0$ e $d_c = 1$ se $c \in A_1$.
21 Essa lista D será usada para percorrer a árvore até a folha de um caractere específico. Essa
22 abordagem utiliza $|H|(1 + o(1))$ bits de memória e tem complexidade de consulta $O(\xi)$, onde
23 ξ é a entropia de ordem zero do texto e é definida como

$$\xi = - \sum_{i=0}^m \frac{f_i}{n} \log \left(\frac{f_i}{n} \right)$$

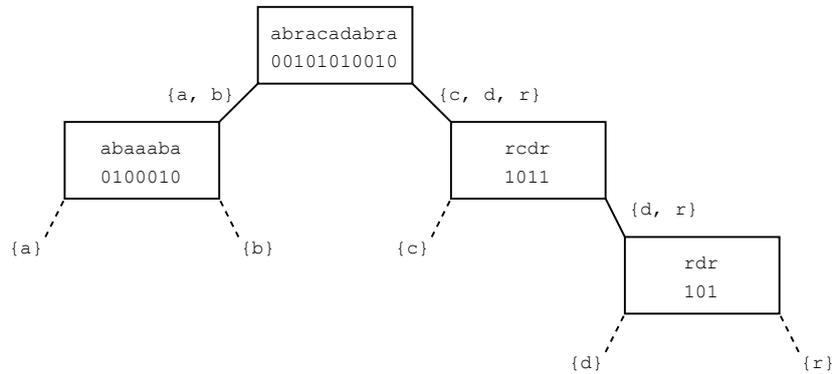


Figura 2.3 A *wavelet tree* balanceada do texto *abracadabra* sobre o alfabeto a, b, c, d, r . Apenas os *bitvectors* são armazenados. O texto e alfabeto são para fins ilustrativos.

1 Um exemplo de uma WT balanceada e uma WT no formato de Huffman podem ser visto
 2 nas Figura 2.3 e Figura 2.4, respectivamente.

3 2.2.1 Construção dinâmica da *wavelet tree* no formato da árvore de Huffman

4 Para a construção da WT no formato da árvore de Huffman de forma dinâmica, modificaremos
 5 a construção descrita em Seção 2.1.1 para não só manter o formato da árvore, como também ar-
 6 mazenar informação em cada nó interno sobre o *bitvector* B e a lista D descritos anteriormente.

7 Suponha que processamos o prefixo $T[0 : i]$. Temos W^i , a WT com formato da árvore de
 8 Huffman para esse prefixo mais um caractere $\#$ que representa todos os caracteres ainda não
 9 vistos. Ao ler o próximo caractere $c = T[i + 1]$, atualizaremos W^i para W^{i+1} . Primeiramente,
 10 adicionaremos o novo *bit* referente à c nos *bitvectors* necessários. Caso c já tenha uma folha
 11 em W^i , percorremos a árvore iniciando da raiz até a folha e, para cada nó interno do caminho,
 12 concatenamos em B um novo *bit* de valor igual à d_c . Caso c não tenha um nó folha em W^i ,
 13 sabemos que, pela construção da árvore de Huffman dinâmica, o nó folha referente a c será
 14 criado vizinho ao nó folha de $\#$. Então iremos percorrer da raiz até a folha de $\#$, concatenando

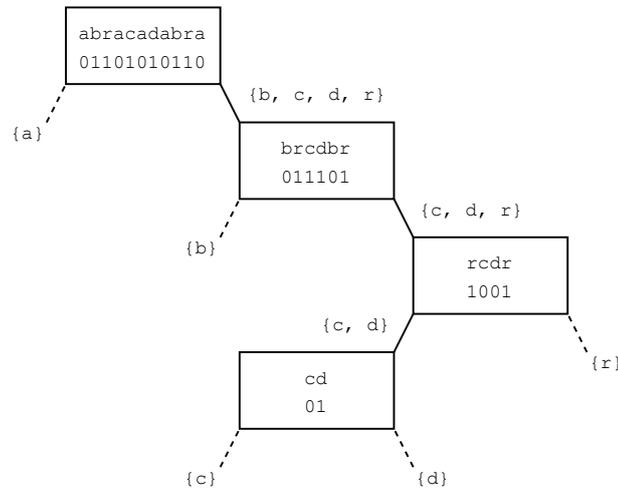


Figura 2.4 A *wavelet tree* no formato Huffman do texto *abracadabra* sobre o alfabeto a, b, c, d, r . Apenas o código de Huffman de cada caractere e os *bitvectors* são armazenados. O texto e alfabeto são para fins ilustrativos.

1 em B um novo *bit* igual à $d_{\#}$ e adicionando $d_c = d_{\#}$. Agora, seguindo a construção da árvore de
 2 Huffman dinâmica, incrementamos em 1 o peso da folha de c caso ela já exista ou criamos os
 3 dois novos nós - um interno e a folha referente à c , e os adicionamos na árvore. Esse novo nó
 4 interno terá, além do peso igual à 1, $B = 0$ e $d_c = 0, d_{\#} = 1$.

5 Após essa modificação, a árvore de Huffman pode sofrer mudanças no seu formato para
 6 manter a propriedade dos irmãos. As mudanças no formato da árvore são realizadas por uma
 7 sequência de trocas entre dois nós e cada troca de posição pode modificar o alfabeto que um
 8 ou mais nós internos representam e, por consequência, modificar tanto B quanto D desses nós.
 9 Para entender quais nós serão modificados em cada troca, é necessário entender o Teorema 2.2.

10 **Teorema 2.2.** Seja $path(x_1, x_2)$ o menor caminho entre dois nós x_1 e x_2 da árvore. Definimos
 11 $LCA(x_1, x_2)$ como sendo o ancestral comum mais profundo (ou *lowest common ancestor*) entre
 12 x_1 e x_2 , ou seja, o nó mais próximo da raiz no caminho entre x_1 e x_2 . Sejam x_i e x_j os nós da
 13 árvore que serão trocados e $x_k = LCA(x_i, x_j)$. O conjunto de nós que precisam ter o bitvector B
 14 e lista D atualizados é igual a $path(x_k, x_i) \cup path(x_k, x_j) \setminus \{x_i, x_j\}$.

1 *Demonstração.* Para cada nó v da árvore, sendo A_i o subconjunto de A que o vértice i repre-
 2 senta, podemos classificá-lo em um dos grupos em relação a um nó $x \in \{x_i, x_j\}$:

3 1. $A_v \subseteq A_x$. Isso significa que v faz parte da subárvore de x . É fácil ver que não precisamos
 4 modificar v nesse caso pois esse nó será movido com x na troca.

5 2. $A_v \cap A_x = \emptyset$. É trivial ver que não precisamos modificar v nesse caso.

6 3. $A_x \subset A_v$. Podemos dividir esse caso em três:

7 $v \in (\text{path}(r, x_k) \setminus \{x_k\})$ Nesse caso v é ancestral de x_k , logo não precisamos mudar os
 8 dados em v porque ambos x_i e x_j continuarão em sua subárvore e continuarão per-
 9 tencendo ao mesmo filho direto.

10 $v = x_k$ Sendo v o $LCA(x_i, x_j)$, precisamos atualizar seus dados pois houve uma inversão
 11 nos *bits* que representam os caracteres em A_{x_i} e A_{x_j}

12 $v \notin (\text{path}(r, x_k) \setminus \{x_k\})$ Neste caso, v é ancestral de x porém ainda descendente de x_k .
 13 Então, como v possui apenas um dos nós trocados como filho, é necessário retirar
 14 as informações do nó que sairá de sua subárvore e adicionar as informações do novo
 15 nó. □

16 Nas Figura 2.5, Figura 2.6 e Figura 2.7 temos exemplos desses três casos para a WT
 17 da palavra *abracadabra*.

18 Sabendo quais nós precisam ser modificados, agora descrevemos como serão modificados.
 19 Sejam x_i e x_j os nós da árvore que serão trocados, $x_k = LCA(x_i, x_j)$ e v um dos nós que será
 20 modificado. Sem perda de generalidade, suponha que v é um ancestral de x_i diferente de x_k .
 21 O bitvector B_v representa as posições de T com caracteres em A_v e a lista D_v que mapeia cada
 22 caractere em A_v para uma das arestas, na direção de sua folha. Por exemplo, na Figura 2.4,
 23 o nó de $T = \text{brcdbr}$ possui $B = 011101$ e $D = \{b : 0, c : 1, d : 1, r : 1\}$. Com a remoção de x_i
 24 da subárvore de v , A_v , B_v e D_v precisam ser modificados para não mais conter os caracteres

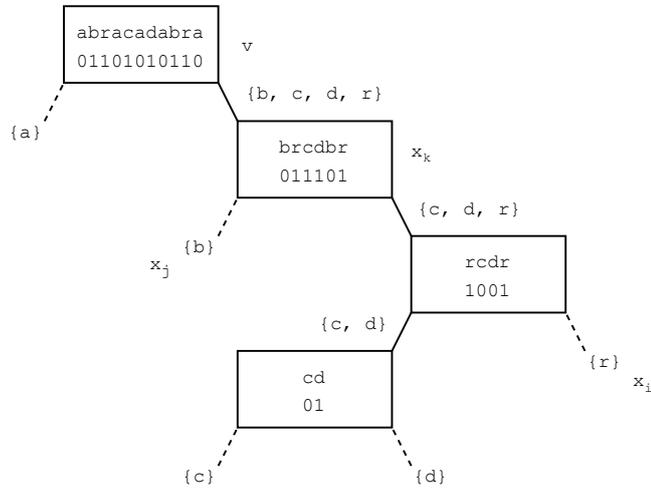


Figura 2.5 Exemplo quando o vértice v não será modificado na troca dos vértices x_i e x_j , onde $LCA(x_i, x_j) = x_k$.

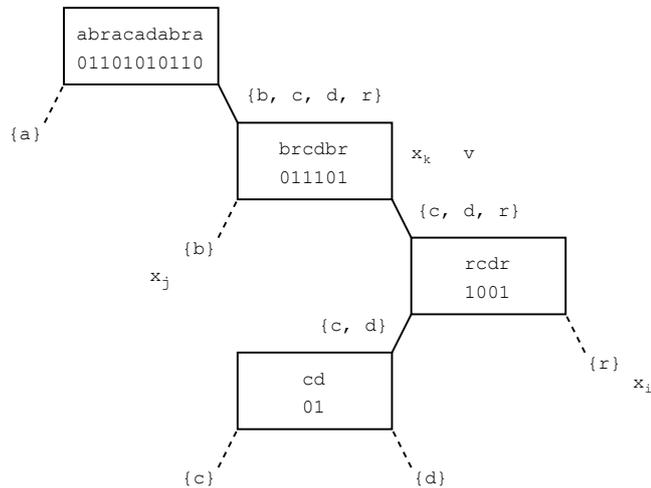


Figura 2.6 Exemplo quando o vértice v será modificado na troca dos vértices x_i e x_j , onde $LCA(x_i, x_j) = x_k$ e $x_k = v$.

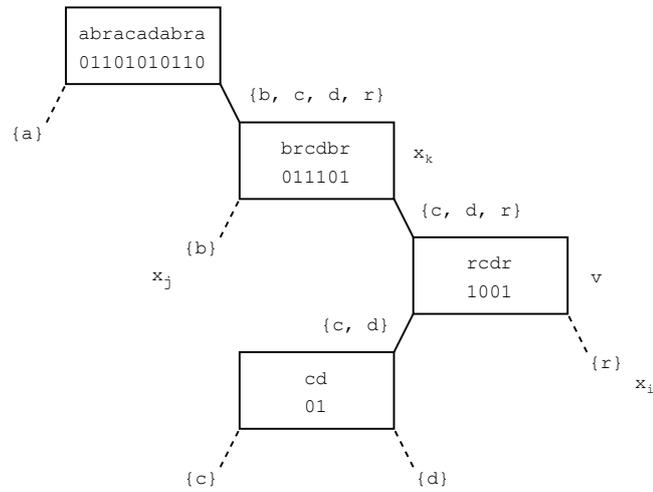


Figura 2.7 Exemplo quando o vértice v será modificado na troca dos vértices x_i e x_j , onde $LCA(x_i, x_j) = x_k$.

1 presentes em A_{x_i} . Com a adição de x_j , é necessário adicionar os caracteres de A_{x_j} em A_v , B_v
 2 e D_v . A adição e remoção em A_v e D_v são triviais de serem realizadas, mas a modificação
 3 de B_v é mais complexa pois não sabemos qual caractere é referente a um *bit* específico, então
 4 precisamos criar uma maneira de achar as posições relevantes para realizar tanto a remoção
 5 quanto a adição.

6 Para o caso que $v = x_k$, não é necessário fazer a operação de remoção ou adição de *bits*, pois
 7 não houve remoção ou adição de nós em sua subárvore. Mas é necessário inverter o *bit* em D_v
 8 e B_v para todo $c \in A_{x_i} \cup A_{x_j}$. Para a inversão de *bits* em B_v , é necessário utilizar os dados das
 9 posições encontrados na remoção dos *bits* dos nós no caminho até x_i e x_j .

2.2.1.1 Remoção de *bits*

11 Para a remoção, começamos de x_i e subimos na árvore até x_k , removendo os *bits* necessários
 12 no caminho. Vamos utilizar a função $select(b, j)$ do *bitvector*, que retorna o índice no vetor do
 13 j -ésimo *bit* b . Por exemplo, para o *bitvector* $B = 110011$ temos

- 1 • $B.select(1,0) = 0$
- 2 • $B.select(0,0) = 2$
- 3 • $B.select(1,3) = 5$

4 A ideia é utilizar a informação de um nó para atualizar o *bitvector* do pai e assim sucessiva-
 5 mente até chegar ao *LCA*. Para atualizar o pai, é necessário saber as posições em seu *bitvector*
 6 que correspondem à caracteres em A_{x_i} . Seja *ant* um nó em $path(x_i, x_k)$, *cur* o pai de *ant* e
 7 P um vetor onde p_i é igual ao índice em B_{cur} do i -ésimo caractere de $proj(T, A_{x_i})$. Inicial-
 8 mente $ant = x_i$, *cur* é o pai de x_i e $P = (p_i = i \mid 0 \leq i < |proj(T, A_{x_i})|)$, que são os índices
 9 em B_{ant} correspondentes aos caracteres de A_{x_i} . Seja y o valor da aresta utilizado para subir
 10 de *ant* para *cur* (0=esquerda, 1=direita). Para modificarmos B_{cur} vamos primeiro atualizar o
 11 vetor P , que está com os índices de B_{ant} , para ter os índices de B_{cur} . Sabendo que em B_{cur}
 12 todo *bit* igual à y se refere à valores em A_{ant} , para atualizarmos os índices em P basta fazer
 13 $(p_i = B_{cur}.select(y, p_i) \mid 0 \leq i < |proj(T, A_{x_i})|)$ e em seguida remover todos esses índices de
 14 B_{cur} . Após a remoção, atualizamos *ant* para ser igual à *cur* e *cur* para ser o novo pai de *ant*.
 15 Esse processo é repetido até que *cur* seja igual à x_k . Quando *cur* for igual à x_k , faremos o
 16 mesmo processo para calcular P mas ao invés de fazermos a remoção dos *bits*, faremos a inver-
 17 são como descrito na demonstração do Teorema 2.2. A figura Figura 2.8 mostra um exemplo
 18 do cálculo de P em uma troca de nós.

19 Quando *cur* for igual a x_k , faremos o mesmo processo para calcular pos mas ao invés de
 20 fazermos a remoção dos *bits*, faremos a inversão.

21 2.2.1.2 Adição de *bits*

22 A adição é realizada após a remoção e será feita do jeito inverso, de cima para baixo, começando
 23 em x_k e descendo até x_j , que será a nova posição de x_i . Vamos utilizar a função $rank(b, j)$ do
 24 *bitvector*, que retorna quantos *bits* iguais a b tem até a posição j . Para o *bitvector* $B = 110011$
 25 temos

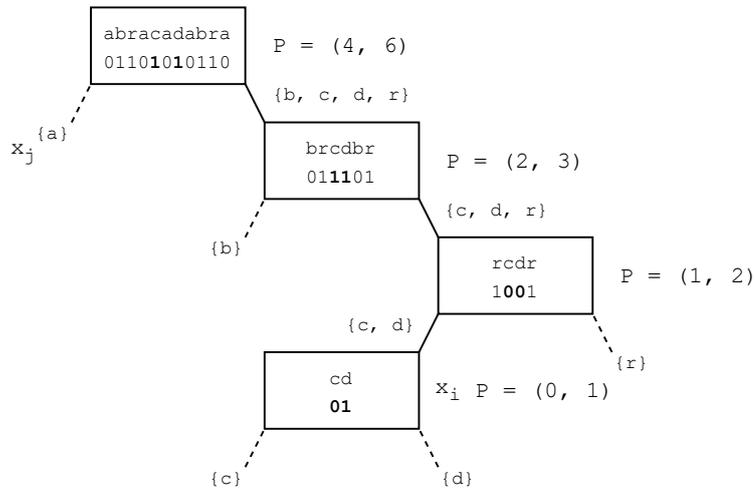


Figura 2.8 Exemplo de valores de P na remoção dos *bits* de x_i na troca com x_j . Os *bits* em negrito serão removidos.

- 1 • $B.rank(1,0) = 1$
- 2 • $B.rank(0,2) = 1$
- 3 • $B.rank(1,5) = 4$
- 4 • $B.rank(0,0) = 0$

5 Seja ant um nó em $path(x_k, x_j)$, cur o filho de ant também em $path(x_k, x_j)$. Seja L um
6 vetor onde l_i é igual posição onde devemos adicionar o i -ésimo *bit* de $proj(T, A_{x_i})$ em B_{cur} .
7 Precisamos processar L para todos os nós no caminho de x_k e x_j e adicionar os *bits* nas posições
8 específicas. Inicialmente, $ant = x_k$ e cur é o filho de x_k mais próximo de x_j . Para o valor inicial
9 de L , usaremos a informação de P calculada na inversão dos *bits* de x_k . Sendo y o valor da
10 aresta utilizado no caminho de ant para cur (0=esquerda, 1=direita), para saber onde inserir
11 os novos *bits* em B_{cur} , basta saber quantos *bits* iguais a y aparecem antes em B_{ant} . Então,
12 $L = (l_i = B_{x_k}.rank(y, p_i) - 1 \mid 0 \leq i < |proj(T, A_{x_i})|)$. Adicionamos em B_{cur} o *bit* y na posição
13 l_i e atualizamos $L = (l_i = B_{cur}.rank(y, l_i) - 1 \mid 0 \leq i < |proj(T, A_{x_i})|)$. Esse processo é repetido
14 até que cur seja igual à x_j , que não precisa ser modificado. A figura Figura 2.9 mostra um

- ¹ mais a pena iterar por todo o *bitvector* ao invés de fazer várias consultas de *rank* e *select*. Além
- ² da *wavelet tree*, é possível retornar o código de Huffman dinâmico do texto.

Resultados

3.1 Dados de teste

Para o teste do método desenvolvido, extraímos textos de Inglês, sequências de DNA e sequências de proteínas de até $3 \cdot 10^6$ caracteres do corpus Pizza&Chili[13].

3.2 Avaliação

Os experimentos foram realizados em um computador com sistema operacional Windows 10, com um processador AMD Ryzen 7 3700X 8-Core 3.60 GHz e 16GB de RAM 3200MHz.

As principais métricas avaliadas serão o tempo de execução para processar todas as atualizações na WT, quantos nós foram trocados e quantos *bits* foram modificados em uma troca de nó em todas atualizações.

Os dados coletados para os diferentes tipos de textos e tamanhos podem ser vistos na Tabela 3.1, Tabela 3.2 e Tabela 3.3. Podemos notar que a quantidade de *bits* trocados ao final de todas atualizações é maior que o tamanho total do texto em todas as tabelas e também na Figura 3.2. Também é importante observar o crescimento linear desse número de trocas ao longo do tamanho do texto, justificando o aumento do tempo de processamento na Figura 3.1. Como discutido na Seção 2.3, o processamento total do texto pode ser quadrático se o número de trocas crescer linearmente ao tamanho do texto, já que depende de ambos fatores. Tamanhos maiores não foram testados por causa do grande tempo de processamento.

Tamanho do texto	Tempo (ms)	Troca de nós	Troca de <i>bits</i>	Alfabeto
10000	80.15	2169	887493	81
50000	236.01	3734	2721971	87
100000	480.87	4515	5396661	90
500000	1958.99	6168	17437093	94
1000000	6020.87	7936	65173480	106
2000000	14985.71	8956	118227207	107
3000000	20987.89	9978	187791347	110

Tabela 3.1 Métricas para textos em inglês de diferentes tamanhos

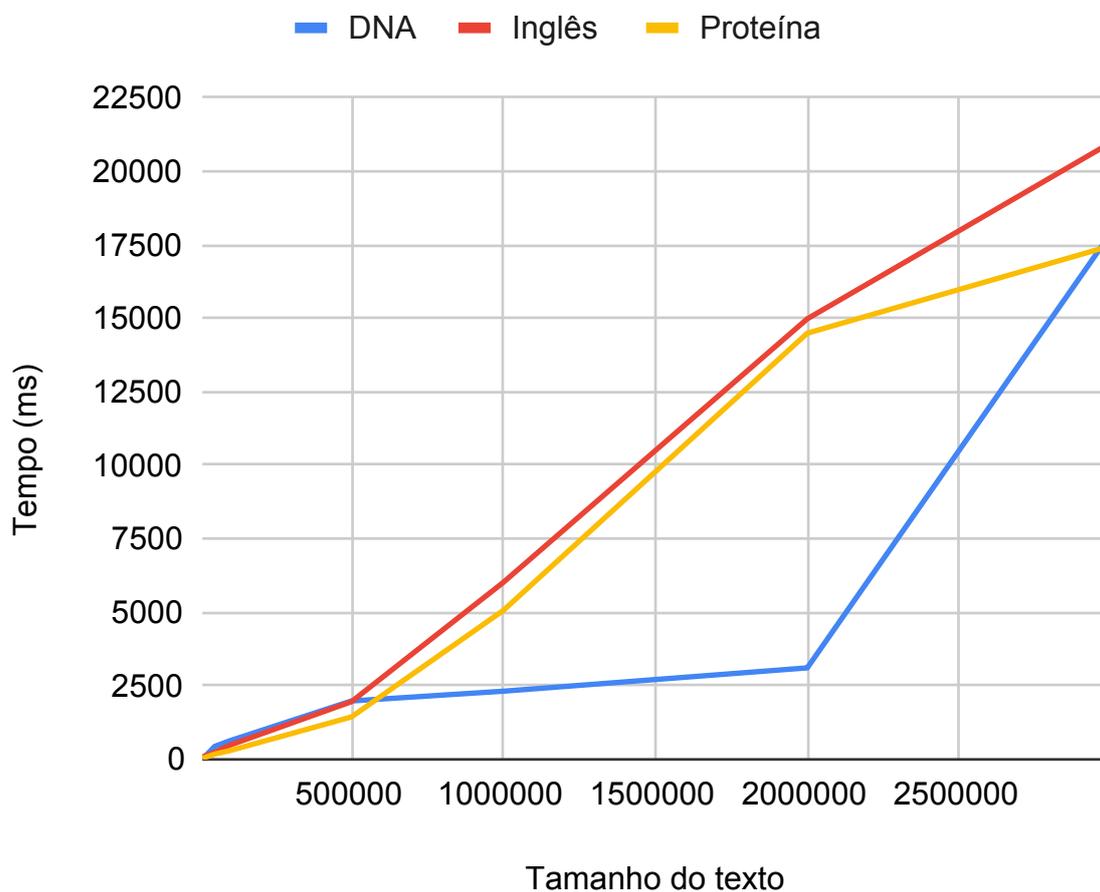


Figura 3.1 Tempo de processamento para cada tamanho de texto de inglês, DNA e proteína.

Tamanho do texto	Tempo (ms)	Troca de nós	Troca de <i>bits</i>	Alfabeto
10000	7.79	72	196331	4
50000	437.06	258	9152145	4
100000	630.80	308	12887742	4
500000	1983.54	393	42864207	4
1000000	2316.17	400	47653798	4
2000000	3109.91	409	62806019	4
3000000	17917.04	553	403465839	4

Tabela 3.2 Métricas para sequências de DNA de diferentes tamanhos

1 Podemos fixar o tamanho do texto e analisar como os diferentes tipos de texto se com-
 2 portam. Nota-se que a sequência de DNA possui um tempo de processamento maior quando
 3 o tamanho é igual à 100000, comparado com os outros dois. Isso acontece porque, como a
 4 sequência de DNA possui um alfabeto menor, a quantidade de *bits* trocados é grande já que os
 5 *bitvectors* estão mais concentrados e possuem uma chance maior de serem um dos nós envol-
 6 vidos na troca, enquanto que os outros dois tipos possuem mais nós com *bitvectors* pequenos.
 7 Esse comportamento fica mais claro quando comparado a coluna de troca de *bits*. Isso pode
 8 ser visto claramente também analisando a diferença de número de troca de *bits* na sequência de
 9 DNA entre os tamanhos 2000000 e 3000000.

10 Analisando os textos de tamanho 3000000, podemos ver que a sequência de proteínas foi
 11 a mais rápida dos três, mas quase igual à de DNA. Olhando apenas para o número de troca
 12 de *bits*, poderia-se estimar que a sequência de DNA deveria ser a mais lenta dos três tipos.
 13 Isso mostra o impacto que um pequeno alfabeto pode ter no tempo de processamento, tendo
 14 quase três vezes mais o número de trocas que a sequência de proteínas e o mesmo tempo de
 15 processamento.

16 Com isso, podemos notar que a complexidade do algoritmo possui termos dependentes
 17 entre si. Mais especificamente, quanto menor o alfabeto, maior as chances de ter uma troca de
 18 nós que afeta muitos *bits*, mas menor será o caminho percorrido até o LCA entre os dois nós da
 19 troca. Isso torna o cálculo do tempo real de processamento do algoritmo algo não trivial para

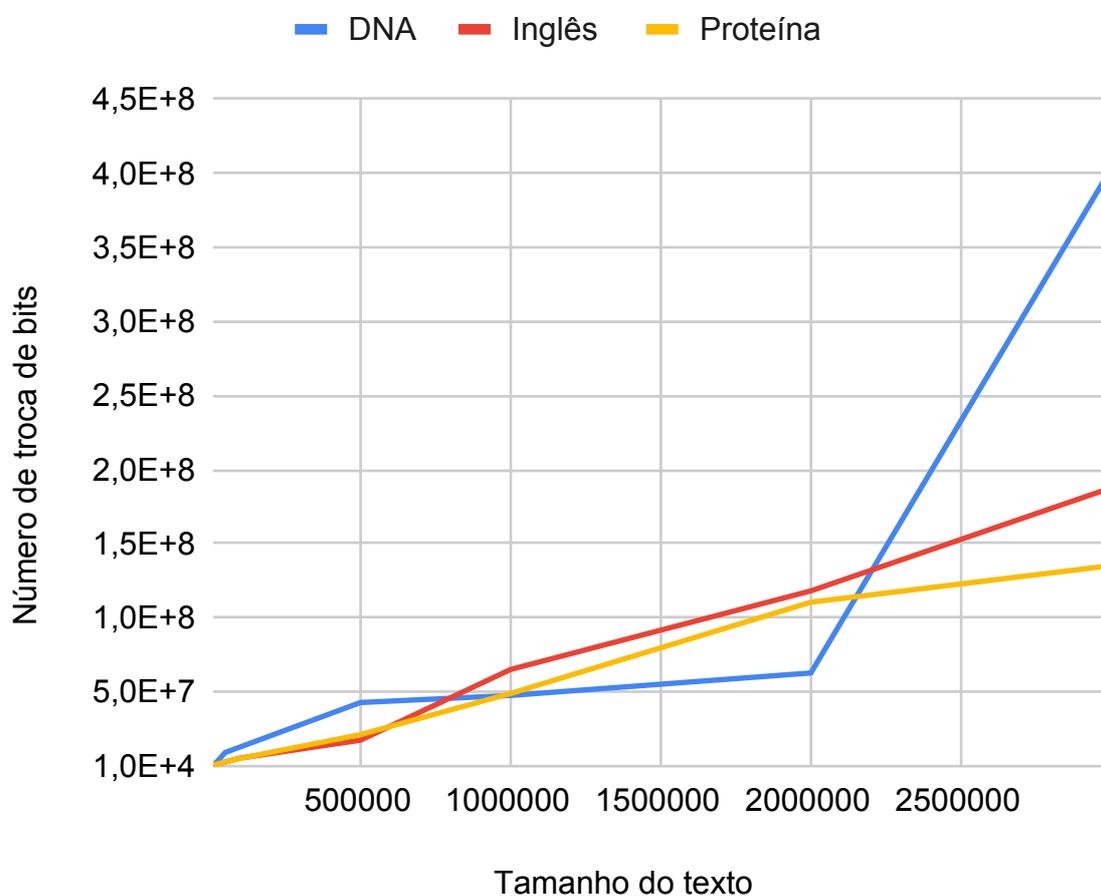


Figura 3.2 Número de troca de *bits* para cada tamanho de texto de inglês, DNA e proteína.

Tamanho do texto	Tempo (ms)	Troca de nós	Troca de <i>bits</i>	Alfabeto
10000	32.70	817	599494	20
50000	175.98	1046	2786206	20
100000	293.82	1146	5182571	20
500000	1439.92	1309	21202639	20
1000000	5066.84	1417	48987888	21
2000000	14497.23	1529	110578342	23
3000000	17472.64	1562	135297287	23

Tabela 3.3 Métricas para sequências de proteínas de diferentes tamanhos

- 1 cada tipo de texto, além de depender da quantidade de nós que será trocada que é diretamente
- 2 relacionada com a entropia.

CAPÍTULO 4

Conclusões

1

2 Nesse trabalho, propusemos um algoritmo para a construção *online* de uma *wavelet tree* no
3 formato da árvore de Huffman. A principal característica do trabalho foi a modificação do
4 algoritmo dinâmico de construção da árvore de Huffman para também armazenar informações
5 necessárias para ser uma *wavelet tree*.

6 A principal vantagem da solução proposta é o fato de não necessitar armazenar o texto em
7 memória, que pode ser uma característica desejada ou necessária. Os resultados experimentais
8 mostram que o tempo de processamento pode ser um problema pelo número de *bits* trocados
9 em mudanças no formato da árvore. Uma possível solução para esse problema seria realizar as
10 mudanças no formato da árvore após processar um bloco do texto ao invés de cada caractere.
11 Isso poderia deixar o formato da árvore mais estável, já que mais mudanças ocorrem quando a
12 diferença de frequência dos caracteres difere em um.

13 Mais melhorias podem ser feitas para alcançar melhor desempenho prático. Por exemplo,
14 conforme discutido na Seção 2.1.1, no Teorema 2.1, dois nós de mesmo peso podem ser consi-
15 derados iguais no quesito de troca de nós e, sabendo disso, podemos fazer uma escolha melhor
16 de nós para se trocar de forma que a quantidade total de *bits* trocados no final seja a menor
17 possível.

18 Outra possível melhoria seria trocar a implementação de *bitvector* utilizada para alguma
19 mais eficiente, já que a atual não possui uma operação de *select* e *rank* otimizada para o *bit* 0.

Referências Bibliográficas

- [1] Maxime Crochemore and Wojciech Rytter. Text Algorithms. Oxford University Press, Inc. 198 Madison Ave. New York, NY, United States, October 1994.
- [2] Paulo G. S. da Fonseca and Israel B. F. da Silva. Online Construction of Wavelet Trees. In Costas S. Iliopoulos, Solon P. Pissis, Simon J. Puglisi, and Rajeev Raman, editors, 16th International Symposium on Experimental Algorithms (SEA 2017), volume 75 of Leibniz International Proceedings in Informatics (LIPIcs), pages 16:1–16:14, Dagstuhl, Germany, 2017. Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik.
- [3] Saska Dönges, Simon J. Puglisi, and Rajeev Raman. On dynamic bitvector implementations. In 2022 Data Compression Conference (DCC), pages 252–261, 2022.
- [4] Newton Faller. An adaptive system for data compression. In Record of the 7-th Asilomar Conference on Circuits, Systems and Computers, pages 593–597, 1973.
- [5] Paolo Ferragina, Raffaele Giancarlo, and Giovanni Manzini. The myriad virtues of wavelet trees. Information and Computation, 207(8):849–866, 2009.
- [6] R. Gallager. Variations on a theme by huffman. IEEE Transactions on Information Theory, 24(6):668–674, 1978.
- [7] Roberto Grossi, Ankur Gupta, and Jeffrey Vitter. High-order entropy-compressed text indexes. Proceedings of the Annual ACM-SIAM Symposium on Discrete Algorithms, 11 2002.

- 1 [8] Roberto Grossi, Ankur Gupta, Jeffrey Scott Vitter, et al. When indexing equals compres-
2 sion: experiments with compressing suffix arrays and applications. In SODA, volume 4,
3 pages 636–645, 2004.
- 4 [9] David A Huffman. A method for the construction of minimum-redundancy codes.
5 Proceedings of the IRE, 40(9):1098–1101, 1952.
- 6 [10] Veli Mäkinen and Gonzalo Navarro. New search algorithms and time/space tradeoffs for
7 succinct suffix arrays. Technical rep. C-2004-20 (April). University of Helsinki, Helsinki,
8 Finland, 2004.
- 9 [11] Veli Mäkinen and Gonzalo Navarro. Succinct suffix arrays based on run-length encoding.
10 In Alberto Apostolico, Maxime Crochemore, and Kunsoo Park, editors, Combinatorial
11 Pattern Matching, pages 45–56, Berlin, Heidelberg, 2005. Springer Berlin Heidelberg.
- 12 [12] Gonzalo Navarro. Wavelet trees for all. Journal of Discrete Algorithms, 25:2–20, 2014.
13 23rd Annual Symposium on Combinatorial Pattern Matching.
- 14 [13] Paolo Ferragina e Gonzalo Navarro. Pizza&chili corpus. [https://pizzachili.](https://pizzachili.dcc.uchile.cl/)
15 [dcc.uchile.cl/](https://pizzachili.dcc.uchile.cl/), Acessado em 2023-01-20.
- 16 [14] Nicola Prezza. A framework of dynamic data structures for string processing. In
17 International Symposium on Experimental Algorithms. Leibniz International Proceedings
18 in Informatics (LIPIcs), 2017.