Genomics

Freie Universität Berlin, Institut für Informatik Knut Reinert, Peter Robinson, Sebastian Bauer Wintersemester 2012/2013 5. Übungsblatt (mit Lösung) vom 24. November 2012 Diskussion am 29. November 2012

1. Name two applications of the Burrows-Wheeler transform. What ist the actual benefit of the Burrows-Wheeler transform in these application (in comparision to other methods)?

Solution See script.

- 2. Given the Burrows-Wheeler transform L = ammnnbsaaaa. (without dot)
 - 1. Decode the original text.
 - 2. Formulate an algorithm that efficiently counts the number of occurrences of a pattern in the original text (without decoding the original text). Describe all of the used data structures.
 - 3. Illustrate how your algorithm works by searching the pattern P =ana.

Solution

(a) We begin with the last character \$ and second last character a in the first rows of *F* and *L*. We iteratively determine preceeding characters using the L-to-F mapping and get the original text T=bananamama\$.

$$LF(i) = C(L[i]) + Occ(L[i], i)$$

$$\Gamma[n-i] = F[LF^{i}(1)]$$

with LF^i denoting the *i*-th function composition with itself.

| a | $\in \Sigma$ | \$ | a | b | n | | | | |
|----|--------------|----|-----|--------------|---|---|--|--|--|
| | n_a | 1 | 5 | 1 | 2 | | | | |
| (| C(a) | 0 | 1 | 6 | 7 | 9 | | | |
| i | F[i] | L | [i] | Occ(L[i], i) | | | | | |
| 1 | \$ | i | a | 1 | | | | | |
| 2 | a | 1 | n | 1 | | | | | |
| 3 | a | 1 | n | 2 | | | | | |
| 4 | a | 1 | n | 1 | | | | | |
| 5 | a | 1 | n | 2 | | | | | |
| 6 | a | 1 | b | 1 | | | | | |
| 7 | b | 5 | \$ | 1 | | | | | |
| 8 | m | i | a | 2 | | | | | |
| 9 | m | i | a | 3 | | | | | |
| 10 | n | i | a | | 4 | | | | |
| 11 | n | i | a | 5 | | | | | |

(b) Algorithm **count** computes the number of occurrences of P[1..m] in T[1..n]:

$$\begin{array}{l} \textit{// count(P[1..m])} \\ i = m, a = 1, b = n; \\ \text{while } ((a < b) \land (i \geq 1)) \text{ do} \\ c = P[i]; \\ a = C(c) + \mathsf{Occ}(c, a - 1) + 1; \\ b = C(c) + \mathsf{Occ}(c, b); \\ i = i - 1; \\ \text{od} \\ \text{if } (b < a) \text{ then return "not found";} \\ else return "found $(b - a + 1) \text{ occurrences";} \end{array}$$$

| (c) | α | \$ | а | b | m | n |
|-----|-------------|----|---|---|---|---|
| | $C[\alpha]$ | 0 | 1 | 6 | 7 | 9 |

- Start with the whole interval $(a_4, b_4) = (1, 11)$.
- Compute L-to-F mappings of the first and last occurrence of a in $L \rightarrow (a_3, b_3) = (2, 6)$.
- Compute L-to-F mappings of the first and last occurrence of n in $L[2..6] \rightarrow (a_2, b_2) = (10, 11)$
- Compute L-to-F mappings of the first and last occurrence of a in $L[10..11] \rightarrow (a_1, b_1) = (5, 6)$
- 3. For the text **tacaacaatacaagag** construct the BWT and the arrays *C* and *Occ*. Use them to search for the pattern **aca**.

| | | | | | | | | | C | $\frac{\in \Sigma}{n_a}$ | \$ 1 0 | a 9 1 | 0 3 1 | <u></u> | g 2 13 | t 2 15 | - | | | |
|----------|------|----|----|----|-----|----|----|----|-------|--------------------------|--------------|-------------|-------------|---------|--------------|--------------|------|--------------|------|--------|
| i | F[i] | | | | | | | | | | | | | | | | L[i] | Occ(L[i], i) | A[i] | Marked |
| 1 | \$ | t | а | С | а | а | С | а | а | t | а | С | а | а | g | а | g | 1 | 17 | * |
| 2 | а | а | С | а | а | t | а | С | а | а | g | а | g | \$ | t | а | С | 1 | 4 | |
| 3 | а | а | g | а | g | \$ | t | а | С | а | а | С | а | а | t | а | С | 2 | 12 | |
| 4 | а | а | t | а | С | а | а | g | а | g | \$ | t | а | С | а | а | С | 3 | 7 | |
| 5 | а | с | а | а | С | а | а | t | а | С | а | а | g | а | g | \$ | t | 1 | 2 | |
| 6 | а | с | а | а | g | а | g | \$ | t | а | С | а | а | С | а | а | t | 2 | 10 | |
| 7 | а | с | а | а | t | а | С | а | а | g | а | g | \$ | t | а | С | а | 1 | 5 | * |
| 8 | а | g | \$ | t | а | С | а | а | С | а | а | t | а | С | а | а | g | 2 | 15 | |
| 9 | а | g | а | g | \$ | t | а | С | а | а | С | а | а | t | а | С | а | 2 | 13 | * |
| 10 | а | t | а | С | а | а | g | а | g | \$ | t | а | С | а | а | С | а | 3 | 8 | |
| 11 | с | а | а | С | а | а | t | а | С | а | а | g | а | g | \$ | t | а | 4 | 3 | |
| 12 | С | а | а | g | а | g | \$ | t | а | С | а | а | С | а | а | t | а | 5 | 11 | |
| 13 | С | а | а | t | а | С | а | а | g | а | g | \$ | t | а | С | а | а | 6 | 6 | |
| 14 | g | \$ | t | а | С | а | а | С | а | а | t | а | С | а | а | g | а | 7 | 16 | |
| 15 | g | a | g | \$ | t | а | С | а | а | С | а | а | t | а | С | a | а | 8 | 14 | |
| 16 | t | a | C | а | а | С | а | а | t | а | С | а | а | g | а | g | \$ | 1 | 1 | * |
| 17 | t | a | С | а | а | g | а | g | \$ | t | а | с | а | ā | с | a | а | 9 | 9 | * |
| . | | | 1 | | - 1 | | ~ | | | | | | | | | | | | | |

First, we count the number of occurrences

- Start with the whole interval $(a_4, b_4) = (1, 17)$.
- Compute L-to-F mappings of the first and last occurrence of a in $L \rightarrow (a_3, b_3) = (2, 10)$.
- Compute L-to-F mappings of the first and last occurrence of c in $L[2...10] \rightarrow (a_2, b_2) = (11, 13)$
- Compute L-to-F mappings of the first and last occurrence of a in $L[11...13] \rightarrow (a_1, b_1) = (5, 7)$

Second, we decode the text starting at each of the interval's rows (5,6, and 7 here). If we decode the text until the beginning, we can conclude the location of the pattern from the number of decoded characters. This, however, would mean in average a running time proportional to the length of the text. For instance, to determine the location of the match defined by row 6, we would need to apply LF 9 times.

If, on the other hand, we would know the position of each character in the text, there would be nothing left to do. Note that the suffix array of the text gives exactly this (denoted as A[i] in the table). Thus, as row 6 is the 6th lexicographical smallest suffix of the text, we can conclude from A[6] = 10 that the pattern must match at position 10.

However, storing the entire suffix array takes additional space. In practice, only a sampled suffix array is stored, meaning that we have to decode the text (i.e., we have to apply LF) until we reach a row with known position. Thus, to achive an $O(m)^1$ running time of the entire match procedure, we need to have access to the suffix array sampled using an O(1) interval.

¹Recall that m is the length of the pattern.

In the table, we specify, which position are known via an asterisk symbol (every 4th position was sampled). Thus, as the A[6] is unknown, we apply LF, which brings us to row 17. As we know the position of that suffix represented by this row and we applied LF only once, the position of our match must be A[17] + 1 = 9 + 1 = 10.