

M1 - BIF

TP3 - Transformée de Burrows Wheeler

Victor Levallois, Claire Lemaitre, Pierre Peterlongo

2023-2024

1 Build BWT

Q 1. Codez (et testez) la fonction $get_BWT(S, SA)$ renvoyant la transformée de Burrows Wheeler BWT du texte S . Pour le calcul du tableau des suffixes, vous pouvez utiliser votre implémentation du TP2 ou utiliser la fonction $SA = \text{simple_kark_sort}(S)$. (from `karkkainan` import `simple_kark_sort`)

2 Reverse BWT

Q 2. Codez la fonction $get_R(BWT)$ renvoyant un tableau tel que $R[i] = \text{rang dans } BWT \text{ du caractère } BWT[i]$.

Q 3. Codez la fonction $get_N(BWT)$ renvoyant un tableau associatif tel que $N[c] = \text{position du premier } c \text{ dans } F$.

Rappel : si $BWT(S)$ est la dernière colonne de la matrice des rotations cycliques de S dans l'ordre lexicographique, F en est la première colonne.

Q 4. Codez la fonction $LF(c, r, N)$ qui renvoie la position dans F du caractère c de rang r .

Q 5. Codez la fonction $BWT2seq(BWT, N, R)$ qui reconstruit, en temps linéaire, la séquence S à partir de sa BWT .

3 Pattern matching - BWT

Q 6. Codez la fonction $find_first(c, i, BWT)$ qui renvoie la position du premier c entre i et $|BWT| - 1$.

Q 7. Codez la fonction $find_last(c, j, BWT)$ qui renvoie la position du dernier c entre 0 et j .

Q 8. A partir des fonctions LF , $find_first$ et $find_last$, codez la fonction $P_in_S(P, BWT, N, R, SA)$ qui renvoie les occurrences de P dans S si elles existent, sinon -1.

4 Bonus

Q 9. Garder les rangs est coûteux en mémoire, implémentez un sous-échantillonnage des rangs pour sacrifier un peu de temps de calcul au profit d'espace mémoire. Pour ce faire, codez `get_R_bis(BWT, p)` qui notera une valeur des rangs de chaque caractère de la BWT tous les p caractères.

Q 10. Idem, pour SA.