

Une chaîne de caractères ASCII est une suite d'octets dont la valeur est comprise entre 0 et 128 (exclue). Un caractère ASCII est donc un octet dont le premier bit est à 0 (zéro). Ce codage ne permet de représenter que les symboles de bases. L'objectif de cet examen est d'écrire les fonctions permettant de gérer les chaînes codées en UTF-8 qui peuvent contenir des symboles accentués ou d'autres alphabets. Un caractère en UTF-8 peut être codé sur 1, 2, 3 ou 4 octets. Le premier octet d'un caractère UTF-8 peut commencer, selon la taille du caractère (en nombre d'octets), soit par 0, 110, 1110 ou 11110. Un octet commençant par 10 correspond à un octet successeur.

- Préfixe 0 0xxx xxxx caractère ASCII sur un octet
- Préfixe 110 110x xxxx 10xxx xxxx caractère sur deux octets, le suivant commençant par 10
- Préfixe 1110 1110 xxxx 10xxx xxxx 10xxx xxxx Trois octets, les deux autres commençant par 10
- Préfixe 11110 1111 0xxx 10xxx xxxx 10xxx xxxx 10xxx xxxx Quatre octets

Exercice 2 : Calcul de la longueur d'une chaîne UTF-8 (2,5 points)

1. Écrire une fonction `int get_prefix(unsigned char c, int n)` qui renvoie le nombre construit à partir des n premiers bits de l'octet c. On attend une fonction d'une ligne avec un simple décalage de bits. Exemple : avec l'octet 215 qui se code 1101 0111, l'appel de `get_prefix(215, 3)` renverra 6 (110 en binaire).

.....

.....

.....

.....

.....

.....

2. Écrire une fonction `int prefix(int n)` qui construit le nombre dont l'écriture en binaire est de longueur n avec n-1 bits 1 suivit d'un unique bit 0. On attend là aussi une fonction d'une ligne avec une simple expression binaire.

- `prefix(1)` doit renvoyer 0 (0 en binaire)
- `prefix(2)` doit renvoyer 2 (10 en binaire)
- `prefix(3)` doit renvoyer 6 (110 en binaire)
- `prefix(4)` doit renvoyer 14 (1110 en binaire)
- `prefix(5)` doit renvoyer 30 (11110 en binaire)

.....

.....

.....

.....

.....

3. En utilisant les deux fonctions précédentes, écrire une fonction `int prefix_p(unsigned char c, int n)` qui renvoie l'équivalent C ANSI du booléen True si l'octet c commence par les bits du préfixe associé à n. Par exemple, `prefix_p(c, 3)` vérifie si c commence par 110 (cf question précédente).

.....

.....

.....

.....

.....

