



Cette page contient la signature de la fonction, ainsi que des fonctions de la même famille, puis une description détaillée, suivi de nombreuses informations diverses (valeur de retour, potentiels codes d'erreur, etc...). Appuyez sur q pour en sortir.

Allez lire la page de manuel de `strcpy`, en particulier la section "BUGS" en bas de la page qui est assez intéressante et montre qu'en informatique, les bugs peuvent causer tout et n'importe quoi : "anything can happen". Autrement dit, d'une exécution à l'autre, un programme buggé ne fera pas la même chose, car le comportement dépendra de l'état précis de la mémoire de la machine, ou du contexte de l'exécution, etc... Il est donc important de programmer de manière à éviter tous les bugs, y compris ceux qui semblent peu probables, et ceux qui semblent inoffensifs.

## La pile, le tas et le segment de données

**Q1.** Le code suivant crée deux chaînes de caractères ayant le même contenu, et affiche leur adresse de stockage :

```
1 #include <stdio.h>
2
3 int main(){
4     char* s = "bonjour";
5     char* t = "bonjour";
6
7     printf("%p\n", s);
8     printf("%p\n", t);
9
10    return 0;
11 }
```

Recopiez et exécutez ce code : que remarquez-vous ?

**Q2.** Essayez maintenant de modifier une des deux chaînes (comme un tableau) en remplaçant une des lettres. Que se passe-t'il ?

En plus de la pile et du tas, une troisième zone, appelée **segment de données**, ou **data segment** en anglais, sert à stocker les variables globales. Par exemple, dans le code suivant :

```
1 int x = 3;
2
3 int main(){
4     ...
5 }
```

la variable `x` est stockée dans le segment de données. De plus, une partie du segment de données est en lecture uniquement, ce qui signifie qu'on ne peut pas modifier son contenu. Ainsi, dans l'exemple précédent avec les deux chaînes de caractères, `s` et `t` sont des variables locales de la fonction `main`, et sont stockées dans la pile, mais elles **pointent vers** la partie lecture uniquement du segment de données, vers la même adresse, où sont stockés les caractères `bonjour`.

Lorsque l'on manipule des chaînes de caractères dans un programme, il faut faire attention à ne pas mélanger les strings alloués sur le tas avec ceux alloués sur la pile ou dans le segment de données.

Q3. On considère le code suivant :

```
1 /* Libère l, liste de n chaînes de caractères. */
2 int free_string_list(char** l, int n){
3     for (int i=0; i<n; i++){
4         free(l[i]);
5     }
6     free(l);
7 }
8
9 int main(){
10     char** l = malloc(3 * sizeof(char*));
11
12     l[0] = "bonjour";
13
14     l[1] = malloc(7*sizeof(char));
15     strcpy(l[1], "coucou");
16
17     l[2] = l[1];
18 }
```

Pour libérer toute la mémoire correctement, faut-il appeler `free(l)` ? `free_string_list(l)` ?  
Aucun des deux ? Et si la liste contenait des milliers de chaînes, allouées parfois dans le tas et parfois pas ?

Dans ce TP, on prendra toujours garde à ce que toute structure de données qui stocke des strings les stocke dans le tas, et à ce qu'elle ne contienne pas de doublons en mémoire : deux strings stockés dans la structure devront pointer vers des **adresses** distinctes. Ainsi, les fonctions de libération de mémoire seront plus simples.

On présente la fonction `strdup` de la librairie `string.h`, qui permet de créer une copie dans le tas d'une chaîne de caractère. Vous pouvez consulter sa documentation dans le manuel pour voir sa spécification précise. Par exemple, le code suivant est une version du code précédent qui utilise `strdup` pour tout réserver dans le tas, il s'exécute sans fuite mémoire :

```
1 /* Libère l, liste de n chaînes de caractères.
2     Chaque case de l doit pointer vers une zone
3     allouée dans le tas, sans doublon */
4 int free_string_list(char** l, int n){
5     for (int i=0; i<n; i++){
6         free(l[i]);
7     }
8     free(l);
9 }
10
11 int main(){
12     char** l = malloc(3 * sizeof(char*));
13
14     l[0] = strdup("bonjour");
15
16     l[1] = malloc(7*sizeof(char));
17     strcpy(l[1], "coucou");
18
19     l[2] = strdup(l[1]);
20     // chaque case de l pointe vers une zone distincte du tas
21     free_string_list(l, 3);
22 }
```

## Dictionnaires et tables de hachage

On rappelle la spécification du type abstrait **dictionnaire** : ce type permet de stocker des associations entre des clés d'un ensemble  $K$  et des valeurs d'un ensemble  $V$ , de telle sorte qu'une clé de  $K$  apparaît au plus une fois dans un dictionnaire. On peut donc voir un dictionnaire comme une fonction partielle de  $K$  dans  $V$ .

- Créer un dictionnaire vide
- Déterminer si une clé  $k$  est dans un dictionnaire  $D$
- Récupérer la valeur associée à une clé  $k$  dans un dictionnaire  $D$
- Supprimer une clé  $k$  d'un dictionnaire  $D$

On rappelle également le principe des tables de hachage : une table de hachage est un tableau de  $m$  cases, munie d'une fonction  $h : K \rightarrow \llbracket 0, m-1 \rrbracket$  appelée fonction de hachage. Les opérations concernant une clé  $k$  se feront dans la case  $i$  du tableau. La fonction de hachage n'étant pas nécessairement injective, deux clés distinctes peuvent être hachées vers le même indice  $i$ , on parle alors de collision.

On propose deux méthodes de résolution des collisions : par chaînage (chaque case du tableau contient une liste chaînée des couples clé-valeur étant entrés en collision dans cette case) et par calcul (lorsqu'une clé est hachée vers un indice de la table contenant déjà une clé distincte, on recalcule un deuxième indice, puis un troisième, etc...).

## Gestion des collisions par chaînage

Le but de cette section est d'implémenter les dictionnaires par table de hachage, en gérant les collisions par chaînage, et de tester l'implémentation sur des exemples simples.

Téléchargez l'archive du TP sur cahier de prépa. Elle contient plusieurs fichiers :

```
keyval.h / keyval.c
dico_chaine.h / dico_chaine.c
dico_hash.h / dico_hash.c
```

Les fichiers `keyval.h/.c` concernent les clés et les valeurs de nos dictionnaires. Y sont implémentées des fonctions pour comparer, afficher et libérer les clés et les valeurs. Ils contiennent également **une fonction de hachage** simple.

Les fichiers `dico_chaine.h/.c` implémentent les dictionnaires par listes chaînées.

Vous n'aurez presque pas à modifier ces fichiers pendant le TP. Lisez les `.h` pour voir l'interface précise proposée, et lisez le `.c` afin de vérifier que le code est cohérent.

Les fichiers `dico_hash.h/.c` sont presque vides, ça sera à vous de les remplir au fil du TP. Ils contiennent pour le moment la définition d'une structure pour les tables de hachage, ainsi qu'un mécanisme permettant de redimensionner la table "à la python" comme vu en cours : lorsque le taux de remplissage passe en dessous de 0.1 ou au dessus de 0.5, on redimensionne la table de façon à ramener le taux de remplissage entre les deux.

- Q4.** Créez un fichier de test, “test.c”, que vous utiliserez pour tester vos fonctions au fur et à mesure que vous avancez dans le TP. A chaque fois que vous écrivez des fonctions, vous devez la tester dans “test.c” avec un bloc de code de la forme :

```
1 int main(){
2     ...
3
4     /* DEBUT TEST nom_fonction_ou_nom_test*/
5     ...
6     ...
7     ...
8     /* FIN TEST nom_fonction_ou_nom_test */
9
10    ...
11    return 0;
12 }
```

Les jeu de tests pourront être des vérifications à l’aide d’assert, ou, lorsque ce n’est pas possible, un simple affichage.

Vous devez régulièrement lancer votre programme de test, et l’analyser avec valgrind pour éviter les fuites mémoires. A la fin du TP, votre programme de test devra avoir des blocs de test pour toutes les fonctions, ne supprimez rien. Pour le moment, en utilisant `assert`, ajoutez des tests permettant de vérifier que `hash(x, m)` renvoie bien  $y$  pour :

- $x = \text{”bonjour”}$ ,  $m = 101$ ,  $y = 53$
- $x = \text{”voici un texte a hacher”}$ ,  $m = 503$ ,  $y = 488$
- $x = \text{”voici un texte a macher”}$ ,  $m = 503$ ,  $y = 236$

- Q5.** Ajoutez un jeu de tests pour la fonction `egal` de “keyval.h”.
- Q6.** Écrivez un jeu de test pour l’implémentation des dictionnaires par listes chaînées, proposée dans `dico_chaine.h/.c`. Votre jeu de test doit couvrir autant d’éventualités que possible (chaîne vide ou presque vide, recherche d’un élément non présent, ajout d’une association pour une clé déjà présente, nombre de clés bien comptabilisé, etc...). Par exemple, créez une chaîne, ajoutez quelques associations dedans, cherchez quelques clés, supprimez quelques associations, recherchez d’autres clés.
- Q7.** Lancez votre jeu de test : vous devez constater qu’il y a quelques erreurs (3 erreurs distinctes à corriger au total). Identifiez et corrigez-les à l’aide de vos tests, puis passez à la suite.

Passons maintenant à l’implémentation des tables de hachage.

- Q8.** Lisez le contenu de `dico_hash.h` pour voir les différentes fonctions proposées dans l’interface. Documentez celles qui ne le sont pas déjà. Regardez également dans `dico_hash.c` pour voir les fonctions d’affichage et de libération qui sont déjà implémentées. Le mécanisme de table redimensionnable est déjà implémenté. Vous pouvez le laisser tel quel, ou bien l’effacer et le recoder par vous-même (si vous êtes très en avance).
- Q9.** Implémentez les fonctions relatives à la structure `hashtable_t` dans `dico_hash.c`, en les testant au fur et à mesure dans `test.c`.
- Q10.** Ajoutez aux dictionnaires l’opération suivante : Renvoyer la liste (sous la forme d’un tableau) des clés d’un dictionnaire.

## A Manipulation de texte

Les dictionnaires ont de nombreuses applications dans l’algorithmique du texte (recherche de motif, compression...). Voyons une application simple plus simple sur le nombre d’occurrences d’un mot.

On considère le texte intégral de “Vingt Mille Lieues sous les mers” de Jules Verne, stocké dans `vingt_mille_lieues.txt`. On veut savoir quel est le mot le plus utilisé parmi ceux faisant plus de 8 lettres. On utilise pour cela l’algorithme suivant :

---

**Algorithme 1** : Mot le plus fréquent

---

**Entrée(s)** :  $t$  un texte,  $K$  un taille de mot minimale

**Sortie(s)** : Le mot de longueur  $\geq K$  ayant le plus d’occurrences dans  $t$

```
1  $D \leftarrow$  dictionnaire vide;
2 pour  $M$  mot de  $t$  faire
3   si  $|M| \geq K$  et  $M \in D$  alors
4      $D[M] \leftarrow D[M] + 1$ ;
5   Else  $D[M] \leftarrow 1$ ;
6 Chercher dans  $D$  la clé  $k_0$  ayant la valeur associée maximale;
7 retourner  $k_0$ 
```

---

- Q11.** Quelle est la complexité de cet algorithme en utilisant une table de hachage avec une bonne fonction de hachage ?
- Q12.** Créez un dossier `texte`, et copiez-collez y les 6 fichiers `.h/.c` de l’implémentation des tables de hachage. Modifiez tout le code nécessaire pour que les clés soient des chaînes de caractères et les valeurs des entiers positifs.
- Q13.** Créez maintenant un fichier “occurrences.c” qui servira à implémenter l’algorithme précédent. Créez-y une fonction `main`, et une fonction `void test()` où vous mettrez les tests.
- Q14.** Écrivez une fonction `KEY* argmax(hashtable_t* d)` renvoyant la clé dont la valeur associée est maximale. On supposera en précondition que le dictionnaire n’est pas vide, et que les valeurs sont des entiers positifs. Écrivez également un jeu de tests afin de tester cette fonction.
- Q15.** Écrivez une fonction `hashtable_t* occurrences(char* filename, int K)` qui renvoie le dictionnaire des occurrences des mots de `filename`, parmi les mots d’au moins  $K$  lettres. N’oubliez pas les questions de gestion mémoire expliquées au début du TP !
- Q16.** Écrivez une fonction `void mot_plus_frequent(char* filename, int K)` qui trouve le mot le plus fréquent dans le fichier `filename`, parmi ceux ayant une longueur au moins  $K$ , et l’affiche ainsi que son nombre d’occurrences. La fonction affichera aussi le nombre total de mots dans le texte. Créez plusieurs petits fichiers permettant de tester cette fonction (pensez aux cas limites).
- Q17.** Une fois que tous vos tests fonctionnent et qu’il n’y a pas de fuite mémoire, vous pouvez commenter la ligne du `main` qui lance les tests. Complétez maintenant le `main` pour obtenir un programme prenant en *argument* (via `argc` et `argv`) : un nom de fichier `fn` et une borne  $K$  et affichant le mot le plus fréquent dans `fn` ayant plus de  $K$  lettres. Testez les performances de votre programme sur `vingt_mille_lieues.txt`.