

# Guide Superflu de programmation en langage C

---

Matthieu Herrb

Version 2.0

Mars 1999

---

Centre National de la Recherche Scientifique  
Laboratoire d'Analyse et d'Architecture des Systèmes

Copyright 1997-1999, Matthieu Herrb. Ce document peut être imprimé et distribué gratuitement dans sa forme originale (comprenant la liste des auteurs). S'il est modifié ou que de des extraits sont utilisés à l'intérieur d'un autre document, alors la liste des auteurs doit inclure tous les auteurs originaux et celui (ou ceux) qui a (qui ont) modifié le document.

Copyright 1997-1999, Matthieu Herrb. This document may be printed and distributed free of charge in its original form (including the list of authors). If it is changed or if parts of it are used within another document, then the author list must include all the original authors AND that author (those authors) who has (have) made the changes.

# Table des matières

<b>I</b>	<b>Quelques pièges du langage C</b>	<b>6</b>
I.1	Fautes de frappe fatales . . . . .	6
I.1.1	Mélange entre = et == . . . . .	7
I.1.2	Tableaux à plusieurs dimensions . . . . .	7
I.1.3	Oubli du <code>break</code> dans les <code>switch</code> . . . . .	8
I.1.4	Passage des paramètres par adresse . . . . .	9
I.2	Problèmes de calcul sur les nombres réels . . . . .	9
I.2.1	Égalité de réels . . . . .	9
I.2.2	Problèmes d'arrondis . . . . .	10
I.2.3	Absence de déclaration des fonctions retournant des doubles . . . . .	10
I.3	Style des déclarations de fonctions . . . . .	11
I.4	Variables non initialisées . . . . .	12
I.5	Ordre d'évaluation indéfini . . . . .	12
I.6	Allocation dynamique de la mémoire . . . . .	13
I.6.1	Référence à une zone mémoire non allouée . . . . .	13
I.6.2	Référence à une zone mémoire libérée . . . . .	13
I.6.3	Libération d'une zone invalide . . . . .	14
I.6.4	Fuites . . . . .	15
I.7	Chaînes de caractères . . . . .	15
I.7.1	Débordement d'une chaîne de caractères . . . . .	15
I.7.2	Écriture dans une chaîne en mémoire statique . . . . .	16
I.8	Pointeurs et tableaux . . . . .	17
I.8.1	Assimilation d'un pointeur et d'un tableau statique . . . . .	17
I.8.2	Appel de <code>free()</code> sur un tableau . . . . .	17
I.9	Entrées/sorties standard . . . . .	17
I.9.1	Contrôle des paramètres de <code>printf</code> et <code>scanf</code> . . . . .	18
I.9.2	Lecture de chaînes de caractères . . . . .	18
I.9.3	Lecture de données binaires . . . . .	19
I.10	Processeurs 64 bits . . . . .	19

I.10.1	Absence de déclarations des fonctions . . . . .	20
I.10.2	Manipulation de pointeurs . . . . .	20
I.11	Pré-processeur . . . . .	20
<b>II</b>	<b>Un peu d'algorithmique</b>	<b>22</b>
II.1	Introduction . . . . .	22
II.2	Allocation dynamique de la mémoire . . . . .	23
II.3	Pointeurs . . . . .	24
II.4	Listes . . . . .	24
II.5	Ensembles . . . . .	24
II.6	Tris et recherches . . . . .	25
II.7	Chaînes de caractères . . . . .	25
<b>III</b>	<b>Créer des programmes sûrs</b>	<b>27</b>
III.1	Comment exploiter les bugs d'un programme . . . . .	28
III.2	Règles pour une programmation sûre . . . . .	30
III.2.1	Éviter les débordements . . . . .	30
III.2.2	Se méfier des données . . . . .	31
III.2.3	Traiter toutes les erreurs . . . . .	32
III.2.4	Limiter les fonctionnalités . . . . .	32
III.2.5	Se méfier des bibliothèques . . . . .	32
III.2.6	Bannir les fonctions dangereuses . . . . .	32
III.3	Pour aller plus loin... . . . . .	33
<b>IV</b>	<b>Questions de style</b>	<b>34</b>
IV.1	Commentaires et documentation . . . . .	35
IV.1.1	Commentaires . . . . .	35
IV.1.2	Documentation . . . . .	37
IV.2	Typologie des noms . . . . .	37
IV.3	Déclarations . . . . .	39
IV.4	Indentation . . . . .	39
IV.5	Boucles . . . . .	40
IV.6	Expressions complexes . . . . .	40
IV.7	Conversion de types . . . . .	41
IV.8	Assertions . . . . .	41
	<b>Références bibliographiques</b>	<b>42</b>
	<b>Index</b>	<b>43</b>

# Introduction

Ce document a pour but de rappeler certaines règles et techniques que tout le monde connaît pour développer une application de taille « sérieuse » en langage C.

Les sources d'informations sur le sujet sont nombreuses, c'est pourquoi je pense que ce document est superflu. En particulier, La Foire Aux Questions (FAQ) du forum Usenet `comp.lang.c` constitue une source d'information beaucoup plus riche, mais en anglais. Elle est accessible par l'URL :

<http://www.eskimo.com/~scs/C-faq/top.html>

Cette FAQ existe aussi sous forme de livre [1].

La bibliographie contient une liste d'autres articles ou ouvrages dont la lecture vous sera plus profitable. Si toutefois vous persistez à vouloir lire ce document, voici un aperçu de ce que vous y trouverez :

Le premier chapitre fait un tour parmi les problèmes les plus souvent rencontrés dans des programmes C.

Le deuxième chapitre donne quelques conseils sur le choix et le type d'algorithmes à utiliser pour éviter de réinventer la roue en permanence.

Le troisième chapitre s'intéresse à l'aspect sécurité informatique des algorithmes et de leurs implémentations, sachant que cet aspect doit être intégré au plus tôt dans l'écriture d'un programme.

Enfin, le quatrième chapitre traite du style du code source : indentation, choix des noms de variables, commentaires,...

Il ne s'agit pas d'obligations à respecter à tout prix, mais suivre ces recommandations permettra de gagner du temps dans la mise au point d'une application et facilitera sa maintenance.

---

# Chapitre I

## Quelques pièges du langage C

---

Le but de ce document n'est pas de faire un cours de langage C. Il y a des livres pour ça. Mais entre les bases du langage et la mise en œuvre concrète de ses fonctionnalités, il y a parfois quelques difficultés.

La plupart des erreurs décrites ici ne sont pas détectées à la compilation.

Certaines de ces erreurs conduisent systématiquement à un plantage du programme en cours d'exécution ou à un résultat faux, alors que d'autres conduisent à des situations que la norme du langage C définit comme « comportements indéterminés », c'est-à-dire que n'importe quoi peut se produire, selon les choix de programmation du compilateur ou du système.

Parfois, dans ce cas, le programme en question à l'air de fonctionner correctement. Une erreur ne se produira que dans certaines conditions, ou bien lorsque l'on tentera de porter le programme en question vers un autre type de machine.

### I.1 Fautes de frappe fatales

Cette section commence par attirer l'attention du lecteur sur quelques erreurs d'inattention qui peuvent être commises lors de la saisie d'un programme et qui ne seront pas détectées par la compilation mais produiront nécessairement des erreurs plus ou moins faciles à détecter à l'exécution.

### I.1.1 Mélange entre = et ==

Cette erreur est l'une des plus fréquentes, elle provient de la syntaxe du langage combinée à l'inattention du programmeur. Elle peut être très difficile à détecter. C'est pourquoi, il est indispensable d'avoir le problème à l'esprit en permanence.

Pour ceux qui ne sauraient pas de quoi il est question ici, rappelez-vous que `x = 0` est une expression valide en C qui affecte à `x` la valeur zéro et qui retourne la valeur affectée, c'est à dire zéro. Il est donc parfaitement légal d'écrire :

```
if (x = 0) {
    /* traitement à l'origine */
    ...
} else {
    /* traitement des autres valeurs */
    ...
}
```

Malheureusement le traitement particulier de  $x = 0$  ne sera jamais appelé, et en plus dans le traitement des  $x \neq 0$  la variable  $x$  vaudra 0!

Pour ceux qui ne l'auraient pas vu, il fallait écrire :

```
if (x == 0) {
    /* traitement à l'origine */
    ...
} else {
    /* traitement des autres valeurs */
    ...
}
```

Certains suggèrent d'utiliser systématiquement la construction suivante, qui a le mérite de provoquer une erreur à la compilation si l'on oublie l'un des « = » :

```
if (0 == x)
```

Cependant, ce problème n'est pas limité au cas de la comparaison avec une constante. Il se pose aussi lorsque l'on veut comparer deux variables.

### I.1.2 Tableaux à plusieurs dimensions

En C les indices d'un tableau à plusieurs dimensions s'écrivent avec autant de paires de crochets qu'il y a d'indices. Par exemple pour une matrice à deux dimensions on écrit :

```
double mat[4][4];
```

```
x = mat[i][j];
```

Le risque d'erreur provient du fait que la notation `mat[i, j]` (qui est employée dans d'autres langages) est également une expression valide en langage C.

### I.1.3 Oubli du `break` dans les `switch`

N'oubliez pas le `break` à la fin de chaque `case` dans une instruction `switch`. Si le `break` est absent, l'exécution se poursuit dans le `case` suivant, ce qui n'est en général pas le comportement voulu. Par exemple :

```
void
print_chiffre(int i)
{
    switch (i) {
        case 1:
            printf("un");
        case 2:
            printf("deux");
        case 3:
            printf("trois");
        case 4:
            printf("quatre");
        case 5:
            printf("cinq");
        case 6:
            printf("six");
        case 7:
            printf("sept");
        case 8:
            printf("huit");
        case 9:
            printf("neuf");
    }
}
```

Dans cette fonction tous les `break` ont été oubliés. Par conséquent, l'appel `print_chiffre(7)` affichera :

```
septhuitneuf
```

Ce qui n'est peut-être pas le résultat escompté.

### I.1.4 Passage des paramètres par adresse

En langage C, les paramètres des fonctions sont toujours passés par valeur : il sont copiés localement dans la fonction. Ainsi, une modification d'un paramètre dans la fonction reste localisée à cette fonction, et la variable de l'appelant n'est pas modifiée.

Pour pouvoir modifier une variable de la fonction appelante, il faut réaliser un passage par adresse explicite. Par exemple, une fonction qui permute deux nombres réels aura le prototype :

```
void swapf(float *x, float *y);
```

Et pour permuter les deux nombres  $x$  et  $y$ , on écrira :

```
swapf(&x, &y);
```

Si on a oublié d'inclure le prototype de la fonction `swap()` avant de l'appeler, le risque est grand d'oublier de passer les adresses des variables. C'est une erreur fréquemment commise avec la fonction `scanf()` et ses variantes.

## I.2 Problèmes de calcul sur les nombres réels

Avant d'attaquer un programme quelconque utilisant des nombres réels, il est indispensable d'avoir pris conscience des problèmes fondamentaux induits par la représentation approchée des nombres réels sur toute machine informatique [2].

Dans de nombreux cas, la prise en compte de ces difficultés se fait simplement, mais il peut s'avérer nécessaire d'avoir recours à des algorithmes relativement lourds, dans le cas par exemple où la précision de la représentation des réels par la machine n'est pas suffisante [3].

### I.2.1 Égalité de réels

Sauf coup de chance, l'égalité parfaite ne peut être obtenue dans le monde réel, il faut donc toujours tester l'égalité à  $\epsilon$  près. Mais il faut faire attention de choisir un  $\epsilon$  qui soit en rapport avec les valeurs à tester.

N'utilisez pas :

```
double a, b;  
...  
if(a == b) /* Faux */
```

Mais quelque-chose du genre :

```
#include <math.h>
```

```
if(fabs(a - b) <= epsilon * fabs(a))
```

qui permet un choix de `epsilon` indépendant de l'ordre de grandeur des valeurs à comparer (À condition que `epsilon` soit strictement positif).

## I.2.2 Problèmes d'arrondis

La bibliothèque standard C propose des fonctions pour convertir des nombres réels en entiers. `floor()` arrondit à l'entier immédiatement inférieur, `ceil()` arrondit à l'entier immédiatement supérieur. Ces fonctions comportent deux pièges :

- elles retournent un type double. Il ne faut pas oublier de convertir explicitement leur valeur en type `int` lorsqu'il n'y a pas de conversion implicite.
- dans le cas d'un argument négatif, elles ne retournent peut-être pas la valeur attendue: `floor(-2.5) == -3` et `ceil(-2.5) == -2`.

La conversion automatique des types réels en entiers retourne quant à elle l'entier immédiatement inférieur en valeur absolue: `(int)-2.3 == -2`.

Pour obtenir un arrondi à l'entier le plus proche on peut utiliser la macro suivante :

```
#define round(x) (int)((x)>0?(x)+0.5:(x)-0.5)
```

## I.2.3 Absence de déclaration des fonctions retournant des doubles

Le type `double` occupe sur la plupart des machines une taille plus importante qu'un `int`. Comme les fonctions dont le type n'est pas déclaré explicitement sont considérées comme retournant un `int`, il y aura problème si la valeur retournée était en réalité un `double` : les octets supplémentaires seront perdus.

Cette remarque vaut pour deux types de fonctions :

- les fonctions système retournant des doubles. L'immense majorité de ces fonctions appartiennent à la bibliothèque mathématique et sont déclarées dans le fichier `math.h`. Une exception à noter est la fonction `strtod()` définie dans `stdlib.h`.
- les fonctions des programmes utilisateur. Normalement toutes les fonctions doivent être déclarées avant d'être utilisées. Mais cette déclaration n'est pas rendue obligatoire par le compilateur. Dans le cas de fonctions retournant un type plus grand qu'un `int`, c'est indispensable. Utilisez des fichiers d'en-tête (`.h`) pour déclarer le type de *toutes* vos fonctions.

## I.3 Style des déclarations de fonctions

L'existence de deux formes différentes pour la déclaration des paramètres des fonctions est source de problèmes difficiles à trouver.

### Style K&R

En C « classique » (également appelé Kernigan et Ritchie ou K&R pour faire plus court), une fonction se déclare sous la forme [4] :

```
int
fonction(a)
    int a;
{
    /* corps de la fonction */
    ...
}
```

Et la seule déclaration possible d'une fonction avant son utilisation est celle du type retourné sous la forme :

```
int fonction();
```

Dans ce cas, tous les paramètres formels de types entiers plus petits que `long int` sont promus en `long int` et tous les paramètres formels de types réels plus petit que `double` sont promus en `double`. Avant l'appel d'une fonction, les conversions suivantes sont effectuées sur les paramètres réels<sup>1</sup> :

- les types entiers (`char`, `short`, `int`) sont convertis en `long int`
- les types réels (`float`) sont convertis en `double`

### Style ANSI

La norme ANSI concernant le langage C a introduit une nouvelle forme de déclaration des fonctions [5] :

```
int
fonction(int a)
{
    /* corps de la fonction */
    ...
}
```

---

1. ici « réel » s'applique à paramètre, en opposition à « formel » et non à « type » (en opposition à « entier »)

Avec la possibilité de déclarer le prototype complet de la fonction sous la forme :

```
int fonction(int a);
```

Si on utilise ce type de déclaration, aucune promotion des paramètres n'est effectuée dans la fonction. De même, si un prototype ANSI de la fonction apparaît avant son appel, les conversions de types effectuées convertiront les paramètres réels vers les types déclarés dans le prototype.

### Mélange des styles

Si aucun prototype ANSI d'une fonction (de la forme `int fonction(int a)`) n'a été vu avant son utilisation, le compilateur peut (selon les options de compilation) effectuer automatiquement les conversions de type citées plus haut, alors qu'une fonction déclarée selon la convention ANSI attend les paramètres avec le type exact qui apparaît dans la déclaration.

Si on mélange les prototypes ANSI et les déclarations de fonctions sous forme K&R, il est très facile de produire des programmes incorrects dès que le type des paramètres est `char`, `short` ou `float`.

## I.4 Variables non initialisées

Les variables déclarées à l'intérieur des fonctions (« automatiques ») sont allouées sur la pile d'exécution du langage et ne sont pas initialisées.

Par contre les variables déclarées statiques sont garanties initialisées à zéro.

## I.5 Ordre d'évaluation indéfini

Sauf exceptions, le C ne définit pas l'ordre d'évaluation des éléments de même précedence dans une expression. Pire que ça, la norme ANSI dit explicitement que le résultat d'une instruction qui dépend de l'ordre d'évaluation n'est pas défini si cet ordre n'est pas défini.

Ainsi, l'effet de l'instruction suivante n'est pas défini : `a[i] = i++`; voici un autre exemple de code dont le comportement n'est pas défini :

```
int i = 3;
printf("%d\n", i++ * i++);
```

Chaque compilateur peut donner n'importe quel résultat, même le plus inattendu dans ces cas. Ce genre de construction doit donc être banni.

Les parenthèses ne permettent pas toujours de forcer un ordre d'évaluation total. Dans ce cas, il faut avoir recours à des variables temporaires.

L'exception la plus importante à cette règle concerne les opérateurs logiques `&&` et `||`. Non seulement l'ordre d'évaluation est garanti, mais en plus l'évaluation est arrêtée dès que l'on a rencontré un élément qui fixe définitivement la valeur de l'expression : faux pour `&&` ou vrai pour `||`.

## I.6 Allocation dynamique de la mémoire

Un des mécanismes les plus riches du langage C est la possibilité d'utiliser des pointeurs qui, combinée avec les fonctions `malloc()` et `free()` ouvre les portes de l'allocation dynamique de la mémoire.

Mais en raison de la puissance d'expression du langage et du peu de vérifications réalisées par le compilateur, de nombreuses erreurs sont possibles.

### I.6.1 Référence à une zone mémoire non allouée

La valeur d'un pointeur désigne l'adresse de la zone mémoire vers laquelle il pointe. Si cette adresse ne correspond pas à une zone de mémoire utilisable par le programme en cours, une erreur (*segmentation fault*) se produit à l'exécution du programme. Mais, même si l'adresse est valide et ne produit pas d'erreur, il faut s'assurer que la valeur du pointeur correspond à une zone allouée correctement (avec `malloc()`, ou sous forme statique par une déclaration de tableau) par le programme.

L'exemple le plus fréquent consiste à référencer le pointeur `NULL`, qui par construction ne pointe vers aucune adresse valable.

Voici un autre exemple de code invalide :

```
int *iptr;
```

```
    *iptr = 1234;
```

```
    printf("valeur : %d\n", *iptr);
```

`iptr` n'est pas initialisé et l'affectation `*iptr = 1234;` ne l'initialise pas mais écrit 1234 à une adresse aléatoire.

### I.6.2 Référence à une zone mémoire libérée

À partir du moment où une zone mémoire a été libérée par `free()`, il est interdit d'adresser son contenu. Si cela se produit, on ne peut pas prédire le comportement du programme.

Cette erreur est fréquente dans quelques cas courants. Le plus classique est la libération des éléments d'une liste chaînée. L'exemple suivant n'est PAS correct :

```
typedef struct LISTE {
    int valeur;
    struct LISTE *suivant;
} LISTE;

void
libliste(LISTE *l)
{
    for (; l != NULL; l = l->suivant) {
        free(l);
    } /* for */
}
```

En effet la boucle `for` exécute `l = l->suivant` *après* la libération de la zone pointée par `l`. Or `l->suivant` référence le contenu de cette zone qui vient d'être libérée.

Une version correcte de `libliste()` est :

```
void
libliste(LISTE *l)
{
    LISTE *suivant;

    for (; l != NULL; l = suivant) {
        suivant = l->next;
        free(l);
    } /* for */
}
```

### I.6.3 Libération d'une zone invalide

L'appel de `free()` avec en argument un pointeur vers une zone non allouée, parce que le pointeur est initialisée vers une telle zone, (cf I.6.1) ou parce que la zone a déjà été libérée (cf I.6.2) est une erreur.

Là aussi le comportement du programme est indéterminé.

### I.6.4 Fuites

On dit qu'il y a fuite de mémoire lorsqu'un bloc alloué par `malloc` n'est plus référencé par aucun pointeur, et qu'il ne peut donc plus être libéré. Par exemple, la fonction suivante, sensée permuter le contenu de deux blocs mémoire, fuit : elle perd le pointeur sur la zone tampon utilisée, sans la libérer.

```
void
mempermute(void *p1, void *p2, size_t length)
{
    void *tmp = malloc(length);
    memcpy(tmp, p1);
    memcpy(p1, p2);
    memcpy(p2, tmp);
}
```

En plus cette fonction ne teste pas le résultat de `malloc()`. Si cette dernière fonction retournait `NULL`, on aurait d'abord une erreur de référence vers une zone invalide.

Pour corriger cette fonction, il suffit d'ajouter `free(tmp)` ; à la fin du code. Mais dans des cas réels, garder la trace des blocs mémoire utilisés, pour pouvoir les libérer n'est pas toujours aussi simple.

## I.7 Chaînes de caractères

Les chaînes de caractères sont gérées par l'intermédiaire des pointeurs vers le type `char`. Une particularité syntaxique permet d'initialiser un pointeur vers une chaîne constante en zone de mémoire statique. Toutes les erreurs liées à l'allocation mémoire dynamique peuvent se produire plus quelques autres :

### I.7.1 Débordement d'une chaîne de caractères

Cela se produit principalement avec les fonctions telles que `gets()`, `strcpy()` ou `sprintf()` qui ne connaissent pas la taille de la zone destination. Si les données à écrire débordent de cette zone, le comportement du programme est indéterminé.

Ces trois fonctions sont à éviter au maximum. La pire de toutes est `gets()` car il n'y a *aucun* moyen d'empêcher l'utilisateur du programme de saisir une chaîne plus longue que la zone allouée en entrée de la fonction.

Il existe 3 fonctions alternatives, à utiliser à la place :

- `fgets()` remplace `gets()`

- `strncpy()` remplace `strcpy()`
- `snprintf()` remplace `sprintf()`. Malheureusement cette fonction n'est pas disponible sur tous les systèmes. Mais il existe un certain nombre d'implémentations « domaine public » de `snprintf()`.

Exemple :

Le programme suivant n'est pas correct :

```
char buf[20];

gets(buf);
if (strcmp(buf, "quit") == 0) {
    exit(0);
}
```

Utilisez plutôt :

```
char buf[20];

fgets(buf, sizeof(buf), stdin);
if (strcmp(buf, "quit") == 0) {
    exit(0);
}
```

## 1.7.2 Écriture dans une chaîne en mémoire statique

La plupart des compilateurs et des éditeurs de liens modernes stockent les chaînes de caractères initialisées lors de la compilation avec des constructions du genre :

```
char *s = "ceci est une chaîne constante\n";
```

dans une zone mémoire non-modifiable. Cela signifie que la fonction suivante (par exemple) provoquera une erreur à l'exécution sur certaines machines :

```
s[0] = toupper(s[0]);
```

L'utilisation du mot-clé `const` permet de détecter cette erreur à la compilation :

```
const char *s = "ceci est une chaîne constante\n";
```

## I.8 Pointeurs et tableaux

Une autre puissance d'expression du langage C provient de la possibilité d'assimiler pointeurs et tableaux dans certains cas, notamment lors du passage des paramètres aux fonctions.

Mais il arrive que cette facilité provoque des erreurs.

### I.8.1 Assimilation d'un pointeur et d'un tableau statique

Il arrive même aux programmeurs expérimentés d'oublier que l'équivalence entre pointeurs et tableaux n'est pas universelle.

Par exemple, il y a une différence importante entre les deux déclarations suivantes :

```
char tableau[] = "ceci est une chaine";  
char *pointeur = "ceci est une chaine";
```

Dans le premier cas, on alloue un seul objet, un tableau de 20 caractères et le symbole `tableau` désigne directement le premier caractère.

Dans le second cas, une variable de type pointeur nommée `pointeur` est allouée d'abord, puis une chaîne constante de 20 caractères et l'adresse de cette chaîne est stockée dans la variable `pointeur`.

### I.8.2 Appel de `free()` sur un tableau

Un tableau est une zone mémoire allouée soit statiquement à la compilation pour les variables globales, soit automatiquement sur la pile pour les variables locales des fonctions. Comme l'accès à ses éléments se fait de manière qui ressemble beaucoup à l'accès aux éléments d'une zone de mémoire allouée dynamiquement avec `malloc()`, on peut les confondre au moment de rendre la mémoire au système et appeler par erreur `free()` avec un tableau en paramètre.

Si les prototypes de `free()` et `malloc()` sont bien inclus dans la portée de la fonction en cours, cette erreur doit au minimum provoquer un warning à la compilation.

## I.9 Entrées/sorties standard

La bibliothèque de gestion des entrées et sorties standard du langage C a été conçue en même temps que les premières versions du langage. Depuis la nécessité de conserver la compatibilité avec les premières versions de cette bibliothèque ont laissé subsister un certain nombre de sources d'erreur potentielles.

### I.9.1 Contrôle des paramètres de `printf` et `scanf`

Les fonctions `printf()` et `scanf()` ainsi que leurs dérivées (`fprintf()`, `fscanf()`, `sprintf()`, `sscanf()`, etc.) acceptent un nombre variable de paramètres de types différents. C'est la chaîne de format qui indique lors de l'exécution le nombre et le type exact des paramètres. Le compilateur ne peut donc pas faire de vérifications. Ainsi, ces fonctions auront un comportement non prévisible si :

- le nombre de paramètres passé est inférieur au nombre de spécifications de conversion (introduites par `%`) dans la chaîne de format,
- le type d'un paramètre ne correspond pas au type indiqué par la spécification de conversion correspondante,
- la taille d'un paramètre est inférieure à la taille attendue par la spécification de conversion correspondante.

Certains compilateurs (dont `gcc`) appliquent un traitement particulier à ces fonctions et vérifient le type des paramètres lorsque le format est une chaîne constante qui peut être analysée à la compilation.

Rappelons les éléments les plus fréquemment rencontrés :

- les paramètres de `scanf()` doivent être les adresses des variables à lire, pas les variables elles-mêmes. Il faut écrire :

```
scanf("%d", &n);
```

et non :

```
scanf("%d", n);
```

- pour lire un double avec `scanf()`, il faut utiliser la spécification de format `%lf`, par contre pour l'afficher avec `printf()` `%f` suffit.

L'explication de cette subtilité se trouve à la section [I.3](#). À vous de la trouver.

### I.9.2 Lecture de chaînes de caractères

La lecture et l'analyse de texte formant des chaînes de caractères est un problème souvent mal résolu. Le cadre théorique général pour réaliser cela correctement est celui de l'analyse lexicographique et syntaxique du texte, et des outils existent pour produire automatiquement les fonctions réalisant ces analyses.

Néanmoins, dans beaucoup de cas on peut se contenter de solutions plus simples en utilisant les fonctions `scanf()`, `fgets()` et `getchar()`.

Malheureusement ces fonctions présentent quelques subtilités qui rendent leur usage problématique.

- `scanf("%s", s)` ; lit un mot de l'entrée standard, séparé par des espaces,

tabulations ou retour à la ligne. Cette fonction saute les séparateurs trouvés à la position courante jusqu'à trouver un mot et s'arrête sur le premier séparateur trouvé après le mot. En particulier si le séparateur est un retour à la ligne, il reste dans le tampon d'entrée.

- `gets(s)` lit une ligne complète, y compris le retour à la ligne final.
- `c = getchar();` et `scanf("%c", &c);` lisent les caractères un à un. La seule différence entre les deux est leur manière de retourner les erreurs en fin de fichier.

Le mélange de ces trois fonctions peut produire des résultats inattendus. En particulier appel à `getchar()` ou `gets()` après

```
scanf("%s", s);
```

retourneront toujours comme premier caractère le séparateur qui a terminé le mot lu par `scanf()`. Si ce séparateur est un retour chariot, `gets()` retournera une ligne vide.

Pour lire des textes comportant des blancs et des retours à la ligne, utilisez exclusivement `fgets()` ou `getchar()`.

L'utilisation de `scanf()` avec le format `%s` est à réserver à la lecture de fichiers structurés simples comportant des mots-clés séparés par des espaces, des tabulations ou des retours à la ligne.

### I.9.3 Lecture de données binaires

Les fonctions `fread()` et `fwrite()` de la bibliothèque des entrées/sorties standard permettent de lire et d'écrire des données binaires directement, sans les coder en caractères affichables.

Les fichiers de données produits ainsi sont plus compacts et plus rapides à lire, mais les risques d'erreur sont importants.

- Le rédacteur et le lecteur doivent être absolument d'accord sur la représentation des types de données de base ainsi écrites.
- Il est fortement conseillé de prévoir une signature du fichier permettant de vérifier qu'il respecte bien le format attendu par l'application qui va le lire.

## I.10 Processeurs 64 bits

De plus en plus de processeurs ont une architecture 64 bits. Cela pose de nombreux problèmes aux utilisateurs du langage C. Beaucoup de programmes ne se compilent plus ou pire, se compile mais ne s'exécutent pas correctement

lorsqu'ils sont portés sur une machine pour laquelle les pointeurs sont plus grands qu'un `int`.

La liste des cas où l'équivalence entiers/pointeurs est utilisée implicitement est malheureusement trop longue et trop complexe pour être citée entièrement. On ne verra que deux exemples parmi les plus fréquents.

### **I.10.1 Absence de déclarations des fonctions**

De même que pour le type double, à partir du moment où les pointeurs n'ont plus la taille d'un entier, toutes les fonctions passant des pointeurs en paramètre ou retournant des pointeurs doivent être déclarées explicitement (selon la norme ANSI) avant d'être utilisées.

Les programmes qui s'étaient contentés de déclarer les fonctions utilisant des double rencontreront des problèmes sérieux.

### **I.10.2 Manipulation de pointeurs**

Il arrive assez fréquemment que des programmes utilisent le type `int` ou `unsigned int` pour stocker des pointeurs avant de les réaffecter à des pointeurs, ou réciproquement qu'ils stockent des entiers dans un type pointeur et cela sans avoir recours à une union.

Dans le cas où les deux types n'ont pas la même taille, on va perdre une partie de la valeur en le transformant en entier, avec tous les problèmes imaginables lorsqu'il s'agira de lui rendre son statut de pointeur.

## **I.11 Pré-processeur**

Le pré-processeur du langage C (`cpp`) pose lui aussi certains problèmes et peut être à l'origine de certaines erreurs.

- certaines erreurs de compilation inexplicables proviennent de la re-définition par le pré-processeur d'un symbole de votre programme. N'utilisez jamais d'identificateurs risquant d'être utilisés aussi par un fichier d'en-tête système dans vos programmes. En particulier, tous les identificateurs commençant par le caractère « souligné » (`_`) sont réservés au système.
- Lors de l'écriture des macros, attention au nombre d'évaluation des paramètres : puisqu'il s'agit de macros et non de fonctions, les paramètres

sont évalués autant de fois qu'ils apparaissent dans le corps de la macro, et non une seule fois au moment de l'appel. Ainsi:

```
#define abs(x) ((x)<0?- (x) : (x))
```

évalue son argument deux fois. Donc `abs(i++)` incrémentera `i` deux fois.

- Comme dans l'exemple précédent, utilisez toujours des parenthèses autour des paramètres dans l'expansion de la macro. C'est le seul moyen de garantir que les différents opérateurs seront évalués dans l'ordre attendu.

---

# Chapitre II

## Un peu d'algorithmique

---

Le but de cette section est donner quelques pistes pour l'utilisation des algorithmes que vous aurez appris par ailleurs, par exemple dans [6].

### II.1 Introduction

Voici en introduction, quelques règles données par R. Pike dans un article sur la programmation en C [7].

La plupart des programmes sont trop compliqués, c'est-à-dire plus compliqués que nécessaire pour résoudre efficacement le problème qui leur est posé. Pourquoi? Essentiellement parce qu'ils sont mal conçus, mais ce n'est pas le but de ce document de discuter de conception, le sujet est trop vaste.

Mais l'implémentation des programmes est également trop souvent compliquée inutilement. Là, les quelques règles suivantes peuvent aider à améliorer les choses.

**Règle 1** On ne peut prédire où un programme va passer son temps. Les goulets d'étranglement se retrouvent à des endroits surprenants. N'essayez pas d'améliorer le code au hasard sans avoir déterminé exactement où est le goulet.

**Règle 2** Mesurez. N'essayez pas d'optimiser un programme sans avoir fait des mesures sérieuses de ses performances. Et refaites-les régulièrement. Si un algorithme plus sophistiqué n'apporte rien, revenez à plus simple.

**Règle 3** Les algorithmes sophistiqués sont lents quand  $n$  est petit, et en général  $n$  est petit. Tant que vous n'êtes pas sûrs que  $n$  sera vraiment grand, n'essayez pas d'être intelligents. (Et même si  $n$  est grand, appliquez d'abord la règle 2).

**Règle 4** Les algorithmes sophistiqués sont plus bugués que les algorithmes simples, parce qu'ils sont plus durs à implémenter. Utilisez des algorithmes et des structures de données simples.

Les structures de données suivantes permettent de traiter tous les problèmes :

- tableaux,
- listes chaînées,
- tables de hachage,
- arbres binaires.

Bien sûr, il peut être nécessaire de les combiner.

**Règle 5** Les données dominent. Si vous avez choisi les bonnes structures de données, les algorithmes deviennent presque évidents. Les structures de données sont bien plus fondamentales que les algorithmes qui les utilisent.

**Règle 6** Ne réinventez pas la roue.

Il existe des bibliothèques de code disponibles librement sur le réseau Internet pour résoudre la plupart des problèmes algorithmiques classiques. *Utilisez-les !*

Il est presque toujours plus coûteux de refaire quelque chose qui existe déjà, plutôt que d'aller le récupérer et de l'adapter.

## II.2 Allocation dynamique de la mémoire

En plus des pièges cités au paragraphe I.6, on peut observer les règles suivantes :

Évitez les allocations dynamiques dans les traitements critiques. L'échec de `malloc()` est très difficile à traiter.

Tenez compte du coût d'une allocation : `malloc()` utilise l'appel système `sbrk()` pour réclamer de la mémoire virtuelle au système. Un appel système est très long à exécuter.

Allouez dynamiquement les objets dont la taille n'est pas connue d'avance et peut varier beaucoup. Il est toujours désagréable d'imposer une taille maximum à un objet parce que le programmeur a préféré utiliser un tableau de taille fixe.

Libérez au plus tôt les objets non utilisés.  
Limitez les copies d'objets alloués dynamiquement. Utilisez les pointeurs.

## II.3 Pointeurs

Les pointeurs sont des outils puissants, même si mal utilisés il peuvent faire de gros dégâts, écrivait Rob Pike dans [7] après s'être planté un ciseau à bois dans le pouce...

Les pointeurs permettent des notations simples pour désigner les objets. Considérons les deux expressions :

```
nodep  
node[i]
```

La première est un pointeur vers un nœud, la seconde désigne un nœud (le même peut-être) dans un tableau. Les deux formes désignent donc la même chose, mais la première est plus simple. Pour comprendre la seconde, il faut évaluer une expression, alors que la première désigne directement un objet.

Ainsi l'usage des pointeurs permet souvent d'écrire de manière plus simple l'accès aux éléments d'une structure complexe. Cela devient évident si l'on veut accéder à un élément de notre nœud :

```
nodep->type  
node[i].type
```

## II.4 Listes

Utilisez de préférence les listes simplement chaînées. Elles sont plus faciles à programmer (donc moins de risque d'erreur) et permettent de faire presque tout ce que l'on peut faire avec des listes doublement chaînées.

Le formalisme du langage LISP est un très bon modèle pour l'expression des opérations sur les listes.

Définissez ou utilisez un formalisme générique pour les listes d'une application.

## II.5 Ensembles

Les ensembles de taille arbitrairement grande sont un peu difficiles à implémenter de manière efficace. Par contre, lorsqu'on a affaire à des ensembles de taille raisonnable (moins d'une centaine d'éléments) et connue d'avance, il est facile de les

implémenter de manière plutôt efficace : l'élément  $n$  de l'ensemble est représenté par le bit  $n \pmod{32}$  de l'élément  $n/32$  d'un tableau d'entiers.

Les opérations élémentaires sur ce type d'ensemble se codent de manière triviale à l'aide des opérateurs binaires  $\&$ ,  $|$ ,  $\sim$ .

## II.6 Tris et recherches

N'essayez pas de programmer un tri. Il existe des algorithmes performants dans la bibliothèque standard C (`qsort()`). Ou bien utilisez les algorithmes de Knuth [8].

Il en est de même pour les problèmes de recherche de données dans un ensemble. Voici un petit éventail des possibilités :

**recherche linéaire** : l'algorithme le plus simple. les éléments n'ont pas besoin d'être triés. La complexité d'une recherche est en  $O(n)$ , si  $n$  est le nombre d'éléments dans la liste. L'ajout d'un élément se fait en temps constant. Cela reste la structure adaptée pour tous les cas où le temps de recherche n'est pas le principal critère. Cette méthode est proposée par la fonction `lsearch()` de la bibliothèque standard C.

**arbres binaires** : les données sont triées et la recherche se fait par dichotomie. Les opérations de recherche et d'ajout se font en  $O(\log(n))$ . Il existe de nombreuses variantes de ce type d'algorithmes, en particulier une version prête à l'emploi fait partie de la bibliothèque C standard : `bsearch()`.

**tables de h-coding** : une fonction de codage (appelée fonction de hachage) associe une clé numérique à chaque élément de l'ensemble des données (ou à un sous-ensemble significativement moins nombreux). Si la fonction de hachage est bien choisie, les ajouts et les recherches se font en temps constant. En pratique, la clé de hachage n'est jamais unique et ne sert qu'à restreindre le domaine de recherche. Une seconde étape faisant appel à une recherche linéaire ou à base d'arbre est nécessaire. Certaines versions de la bibliothèque standard C proposent la fonction `hsearch()`.

## II.7 Chaînes de caractères

Les chaînes de caractères donnent lieu à de nombreux traitements et posent pas mal de problèmes algorithmiques. Voici quelques conseils pour une utilisation

saine des chaînes de caractères :

- évitez de limiter arbitrairement la longueur des chaînes : prévoyez l'allocation dynamique de la mémoire en fonction de la longueur.
- si vous devez limiter la longueur d'une chaîne, vérifiez qu'il n'y a pas débordement et prévoyez un traitement de l'erreur.
- utilisez de préférence les fonctions de lecture caractère par caractère pour les chaînes. Elles permettent les meilleures reprises en cas d'erreur.
- prévoyez à l'avance l'internationalisation de votre programme : au minimum, considérez que l'ensemble des caractères à traiter est celui du codage ISO Latin-1.
- utilisez les outils `lex` et `yacc` pour les traitements lexicographiques et syntaxiques un peu complexes : vos programmes gagneront en robustesse et en efficacité.

---

# Chapitre III

## Créer des programmes sûrs

---

Bien souvent un programmeur se satisfait d'un programme qui a l'air de fonctionner correctement parce que, apparemment, il donne un résultat correct sur quelques données de test en entrée. Que des données complètement erronées en entrée produisent des comportements anormaux du programme ne choque pas outre-mesure.

Certaines catégories de programmes ne peuvent pas se contenter de ce niveau (peu élevé) de robustesse. Le cas le plus fréquent est celui de programmes offrant des services à un grand nombre d'utilisateurs potentiels, sur le réseau Internet par exemple, auxquels on ne peut pas faire confiance pour soumettre des données sensées en entrée. Cela est particulièrement crucial pour les programmes s'exécutant avec des privilèges particuliers (par exemple exécutés sous l'identité du super-utilisateur sur une machine Unix).

En effet, les bugs causés par les débordements de tableaux ou les autres cas d'écrasement de données involontaires peuvent être utilisés pour faire exécuter à un programme du code autre que celui prévu par le programmeur. Lorsque ce code est le fruit du hasard, (des données brusquement interprétées comme du code), l'exécution ne vas pas très loin et se termine généralement par une erreur de type « bus error » ou « segmentation violation ».

Par contre, un programmeur mal intentionné peut utiliser ces défauts en construisant des jeux de données d'entrée qui font que le code exécuté accidentellement ne sera plus réellement le fruit du hasard, mais bel et bien un morceaux de

programme préparé intentionnellement et destiné en général à nuire au système ainsi attaqué [9].

Par conséquent, les programmes ouverts à l'utilisation par le plus grand nombre doivent être extrêmement vigilants avec toutes les données qu'ils manipulent.

De plus, comme il est toujours plus facile de respecter les règles en les appliquant dès le début, on gagnera toujours à prendre en compte cet aspect sécurité dans tous les programmes, même si initialement ils ne semblent pas promis à une utilisation sensible du point de vue sécurité.

Garfinkel et Spafford ont consacré le chapitre 22 de leur livre sur la sécurité Unix et internet [10] à l'écriture de programme sûrs. Leurs recommandations sont souvent reprises par d'autres auteurs.

La robustesse supplémentaire acquise par un programme qui respecte les règles énoncées ici sera toujours un bénéfice pour l'application finale, même si les aspects sécurité ne faisaient pas partie du cahier des charges initiales. Un programme conçu pour la sécurité est en général aussi plus robuste face aux erreurs communes, dépourvues d'arrière pensées malveillantes.

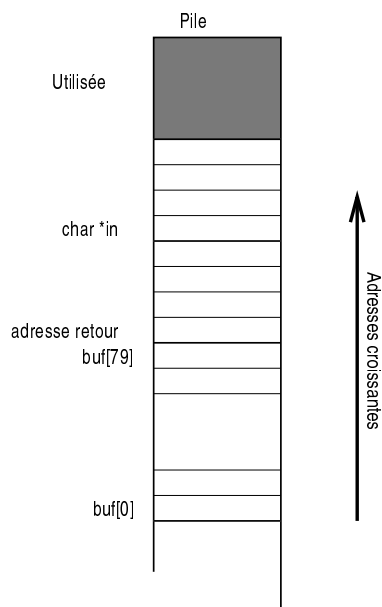
### III.1 Comment exploiter les bugs d'un programme

Un exposé exhaustif des techniques utilisées par les « pirates » informatiques pour exploiter les bugs ou les erreurs de conception d'un logiciel dépasse largement le cadre de ce document. Nous allons simplement présenter ici un exemple classique de bug qui était présent dans des dizaines (voire des centaines) de programmes avant que l'on ne découvre la manière de l'exploiter pour faire exécuter un code arbitraire au programme qui le contient.

Voici une fonction d'un programme qui recopie la chaîne de caractères qu'elle reçoit en argument dans un buffer local (qui est situé sur la pile d'exécution du programme).

```
int
maFonction(char *in)
{
    char buf[80];

    strcpy(buf, in);
    ...
    return 0;
}
```

FIG. III.1 – Organisation des données sur la pile lors dans *maFonction*.

Lors de l'exécution de cette fonction, l'organisation des données sur la pile sera celle décrite sur la figure III.1.

Sur cette figure, il apparaît clairement qu'un débordement par le haut du tableau `buf` va écraser sur la pile l'adresse de retour de la fonction. Dans le cas d'une erreur involontaire, cela conduira à un saut à une adresse invalide et provoquera donc une erreur du type *segmentation violation*.

Par contre, on peut exploiter cette lacune pour faire exécuter au programme en question un morceau de code arbitraire. Il suffit pour cela de s'arranger pour que les quatre premiers octets du débordement soient une adresse donnée sur la pile, dans la zone déjà allouée, et que le reste du débordement soit un programme en code machine de la bonne longueur pour commencer pile à l'adresse que l'on a mise dans l'adresse de retour.

Ainsi, à la fin de l'exécution de `maFonction`, le processeur va dépiler une mauvaise adresse de retour et continuer son exécution dans le code passé en excès dans le tableau `in`.

Ainsi exposée, cette technique semble assez rudimentaire et difficile à mettre en œuvre. Il est cependant courant de trouver sur Internet dans des forums spécialisés des scripts tout faits capables d'exploiter ce type de vulnérabilité dans les applications les plus courantes des systèmes existants.

## III.2 Règles pour une programmation sûre

La liste des règles qui suivent n'est pas impérative. Des programmes peuvent être sûrs sans respecter ces règles. Elle n'est pas non plus suffisante : d'une part parce qu'il est impossible de faire une liste exhaustive (on découvre chaque semaine de nouvelles manières d'exploiter des programmes apparemment innocents), et d'autre part parce que seule une conception rigoureuse permet de protéger un programme contre les erreurs volontaires ou non des utilisateurs.

### III.2.1 Éviter les débordements

C'est la règle principale. Il ne faut jamais laisser la possibilité à une fonction d'écrire des données en dehors de la zone mémoire qui lui est destinée. Cela peut paraître trivial, mais c'est cependant des problèmes de ce type qui sont utilisés dans la grande majorité des problèmes de sécurité connus sur Internet.

Il y a en gros deux techniques pour arriver à cela. Je ne prendrai pas partie pour une technique ou une autre, par contre il est relativement évident que l'on ne peut pas (pour une fois) les mélanger avec succès.

- **Allouer dynamiquement toutes les données.** En n'imposant aucune limite statique à la taille des données, le programme s'adapte à la taille réelle des données et peut toujours les traiter sans risque d'erreur, à condition que la machine dispose de suffisamment de mémoire.

C'est d'ailleurs là que réside la difficulté majeure de cette méthode. Lorsque la mémoire vient à manquer, il est souvent très délicat de récupérer l'erreur proprement pour émettre un diagnostic clair, libérer la mémoire déjà allouée mais inutilisable et retourner dans un état stable pour continuer à fonctionner.

Une autre difficulté provient de la difficulté dans certains cas de prédire la taille d'une donnée. On se trouve alors contraint de réallouer plusieurs fois la zone mémoire où elle est stockée, en provoquant autant de recopies de ces données, ce qui ralentit l'exécution du programme.

- **Travailler uniquement avec des données de taille fixe allouées statiquement.** Avec cette technique on s'interdit tout recours à la fonction `malloc()` ou à ses équivalents. Toutes les données extérieures sont soit tronquées pour contenir dans les buffers de l'application soit traitées séquentiellement par morceaux suffisamment petits. Dans ce cas le traitement des erreurs est plus simple, par contre certains algorithmes deviennent complexes. Par exemple, comment réaliser par exemple un tri de données arbitrairement grandes sans allocation dynamique de mémoire?

Il faut remarquer ici que les systèmes de mémoire virtuelle aident à gommer les défauts respectifs des deux approches. La possibilité d'allouer des quantités de mémoire bien supérieures à la mémoire physique disponible aide à retarder l'apparition du manque de mémoire dans le premier cas. La faculté de la mémoire virtuelle à ne réclamer de la mémoire physique que pour les données réellement référencées permet dans la seconde approche de prévoir des tailles de données statiques relativement grandes sans monopoliser trop de ressources si les données réelles sont très souvent beaucoup plus petites.

### III.2.2 Se méfier des données

Si vous considérez que l'utilisateur de votre programme veut nuire à votre système, toutes les données qu'il fournit en entrée sont potentiellement dangereuses. Votre programme doit donc analyser finement ses entrées pour rejeter intelligemment les données suspectes, sans pénaliser outre mesure les utilisateurs honnêtes.

Le premier point à vérifier a déjà été évoqué : il faut éviter que la taille des données d'entrée ne provoque un débordement interne de la mémoire. Mais il y a d'autres points à vérifier :

- Vérifier les noms des fichiers. En effet l'utilisateur peut essayer d'utiliser les privilèges potentiels de votre programme pour écraser ou effacer des fichiers système.
- Vérifier la syntaxe des commandes. Chaque fois qu'un programme commande l'exécution d'un script, il est possible d'exploiter la syntaxe particulièrement riche du shell Unix pour faire faire à une commande autre chose que ce pourquoi elle est conçue. Par exemple si un programme exécute le code suivant pour renommer un fichier :

```
printf(cmd, "mv %s %s.bak", fichier, fichier);
system(cmd);
```

pour renommer un fichier, si `fichier` contient la valeur :

```
toto toto.bak ; cat /etc/passwd ;
```

la fonction `system()` va exécuter :

```
mv toto toto.bak ; cat /etc/passwd ; toto toto.bak ; cat /etc/passwd;
```

Ce qui va afficher le fichier des mots de passe cryptés du système en plus du résultat initialement attendu.

- Vérifier l'identité de l'utilisateur. Dans tous les cas où cela est possible, l'identification des utilisateurs ne limite pas directement ce qu'il peuvent faire, mais aide à mettre en place des mécanismes de contrôle d'accès.

### III.2.3 Traiter toutes les erreurs

Éviter que des erreurs se produisent n'est pas toujours possible.

Prévoyez des traces des problèmes de sécurité, non visibles directement par l'utilisateur. Par contre il est indispensable de prévenir l'utilisateur de l'existence et de la nature de ces traces (loi Informatique et libertés + effet dissuasif)

### III.2.4 Limiter les fonctionnalités

Certaines fonctionnalités d'un programme peuvent être dangereuses. Il faut y songer dès la spécification pour éviter de fournir au pirate les moyens de parvenir facilement à ses fins.

- possibilité de créer un shell
- affichage de trop d'information
- traitements sans limites

### III.2.5 Se méfier des bibliothèques

Les règles ci-dessus devraient être respectées par les programmeurs qui ont réalisé les bibliothèques de fonctions utilisées par votre application (bibliothèque C standard, interface graphique, calcul matriciel,...). Mais en êtes-vous certains? L'expérience montre que des problèmes du style de ceux évoqués ici sont présents dans de nombreux logiciels commerciaux, comme dans les logiciels libres.

En général il n'est pas possible d'auditer tout le code des bibliothèques utilisées, soit parce que celui-ci n'est pas disponible, soit simplement par manque de temps et/ou de compétence.

Interrogez-vous sur le type d'algorithme utilisé par les fonctions appelées par votre code. Si l'un d'eux présente des risques de débordement interne, vérifiez deux fois les données, à l'entrée et à la sortie pour détecter du mieux possible les éventuelles tentatives d'attaque. En cas de doute sur un résultat votre programme doit le signaler au plus tôt, et ne pas utiliser ce résultat.

### III.2.6 Bannir les fonctions dangereuses

Certaines fonctions de la bibliothèque C standard sont intrinsèquement dangereuses parce que leur sémantique ne permet pas de respecter les règles présentées ci-dessus. Il faut donc s'interdire impérativement de les utiliser. Il peut y avoir des cas où ces fonctions peuvent être utilisées malgré tout de manière sûre. A mon avis, même dans ces cas, il faut les éviter et leur préférer une version sûre. Cela

facilite la vérification a posteriori du code, en évitant de provoquer des fausses alarmes qui peuvent être coûteuses à désamorcer. De plus, le raisonnement qui vous a amené à considérer une utilisation d'une fonction dangereuse comme sûre peut être faux ou incomplet et donc le risque n'est pas éliminé complètement.

Ne pas utiliser	remplacer par	remarque(s)
gets()	fgets()	risque de débordement
scanf()	strtol(), strtod(), strtok(),...	idem
sprintf()	snprintf()	risque de débordement
strcat()	strncat()	idem
strcpy()	strncpy()	idem
mktemp()	mkstemp()	section critique: entre la création et l'ouverture du fichier temporaire
system()	fork() & exec()	possibilité d'exploiter le shell

### III.3 Pour aller plus loin...

Le numéro d'avril 1998 du magazine électronique *Sunworld Online* propose une méthodologie de développement de logiciels sûrs (SDSDM – Software Development Security Design Methodology). L'article est accessible à l'adresse: <http://www.sunworld.com/swol-04-1998/swol-04-security.html>

Adam Shostack, consultant en sécurité informatique, a rédigé un guide pour la relecture du code destiné à tourner dans un firewall: <http://www.homeport.org/~adam/review.html> qui est souvent cité comme référence pour l'évaluation des logiciels de sécurité.

Le projet FreeBSD propose un ensemble de recommandations au développeurs qui souhaitent contribuer au projet : <http://www.freebsd.org/security/programmers.html>

Matt Bishop a été l'un des précurseurs de la notion de programmation robuste. Ses articles sur le sujet font référence. <http://olympus.cs.ucdavis.edu/~bishop/secprog.html>

Enfin, le chapitre de [10] consacré à la programmation sûre est disponible en ligne: [ftp://ftp.auscert.org.au/pub/auscert/papers/secure\\_programming\\_checklist](ftp://ftp.auscert.org.au/pub/auscert/papers/secure_programming_checklist)

---

# Chapitre IV

## Questions de style

---

Les règles présentées ici ne sont pas impératives, il s'agit juste d'exemples de bonnes habitudes qui facilitent la relecture d'un programme.

Ce qui est le plus important, c'est de penser qu'un programme doit pouvoir être lu et compris par quelqu'un d'extérieur, qui ne connaît pas forcément tout du logiciel dont est extrait le morceau qu'il relit. La lisibilité d'un code source (qui peut s'analyser avec les règles de la typographie) est une très bonne mesure de sa qualité.

Les guides de style pour les programmeurs C sont très nombreux dans la littérature. Presque chaque ouvrage sur le langage propose son style. Toutes les grandes entreprises et les grands projets de logiciel ont leurs règles.

Un guide a servi de modèle à de nombreux programmeurs : le *Indian Hill C Style and Coding Standards* des Laboratoires Bells [11]. Ce guide a été amendé et modifié de très nombreuses fois, mais sert de référence implicite commune à de nombreux autres guides.

## IV.1 Commentaires et documentation

### IV.1.1 Commentaires

Bien commenter un programme est sans doute la chose la plus difficile de toute la chaîne de développement. C'est un des domaines où « le mieux est l'ennemi du bien » s'applique avec le plus d'évidence.

De manière générale, le commentaire permet d'apporter au lecteur d'un programme une information que le programme lui-même ne fournit pas assez clairement.

Pour évaluer la qualité d'un commentaire, le recours aux règles de la typographie est précieux : un commentaire ne doit pas être surchargé de ponctuation ou de décorations. Plus il sera sobre, plus il sera lisible.

Un bon commentaire a surtout un rôle introductif : il présente ce qui suit, l'algorithme utilisé ou les raisons d'un choix de codage qui peut paraître surprenant.

Un commentaire qui paraphrase le code et vient après coup n'apporte rien. De même, il vaut mieux un algorithme bien programmé et bien présenté avec des noms de variables bien choisis pour aider à sa compréhension, plutôt qu'un code brouillon très compact suivi ou précédé de cent lignes d'explications. L'exemple extrême du commentaire inutile est :

```
i++; /* ajoute 1 a la variable i */
```

D'ailleurs, avec ce genre de commentaires, le risque de voir un jour le code et le commentaire se contredire augmente considérablement.

Enfin, il est sûrement utile de rappeler *quand* écrire les commentaires : tout de suite en écrivant le programme. Prétendre repasser plus tard pour commenter un programme c'est une promesse d'ivrogne qui est très difficile à tenir.

### En-têtes de fichiers

Il peut être utile d'avoir en tête de chaque fichier source un commentaire qui attribue le copyright du contenu à son propriétaire, ainsi qu'un cartouche indiquant le numéro de version du fichier, le nom de l'auteur et la date de la dernière mise à jour, avant une description rapide du contenu du fichier.

Exemple (ici l'en-tête est maintenue automatiquement par RCS) :

```
/**
*** Copyright (c) 1997,1998 CNRS-LAAS
***
*** $Source: /home/matthieu/cvs/doc/cours/C/style.tex,v $
*** $Revision: 1.9 $
```

```

*** $Date: 1999/04/05 19:45:49 $
*** $Author: matthieu $
***
*** Fonctions de test sur les types du langage
***/

```

Remarque: la notice de copyright rédigée ainsi n'a aucune valeur légale en France. Pour une protection efficace d'un programme, il faut le déposer auprès d'un organisme spécialisé. Néanmoins, en cas de conflit, la présence de cette notice peut constituer un élément de preuve de l'origine du logiciel.

### En-têtes de fonctions

Avant chaque fonction, il est très utile d'avoir un commentaire qui rappelle le rôle de la fonction et de ses arguments. Il est également très utile d'indiquer les différentes erreurs qui peuvent être détectées et retournées.

Exemple :

```

/**
** insertAfter - insère un bloc dans la liste après un bloc
**                particulier
**
** paramètres:
**  list: la liste dans laquelle insérer le bloc. NULL crée une
**        nouvelle liste
**  pBloc: pointeur sur le bloc a insérer
**
** retourne: la nouvelle liste
**/

```

Donner le type des paramètres ne sert à rien, car la déclaration formelle de la fonction, avec le type exact suit immédiatement.

### Commentaires dans le code

On peut presque s'en passer si l'algorithme est bien présenté (voir aussi remarque sur la complexité au chapitre II).

Il vaut mieux privilégier les commentaires qui répondent à la question *pourquoi?* par rapport à ceux qui répondent à la question *comment?*

Les commentaires courts peuvent être placés en fin de ligne. Dès qu'un commentaire est un peu long, il vaut mieux faire un *bloc* avant le code commenté.

Pour un bloc, utiliser une présentation sobre du genre :

```
/*  
 * Un commentaire bloc  
 * Le texte est mis en page de manière simple et claire.  
*/
```

Les commentaires placés dans le code doivent être indentés comme le code qu'ils précèdent.

N'essayez pas de créer des cadres compliqués, justifiés à droite ou avec une apparence 3D. Cela n'apporte aucune information, et est très dur à maintenir propre lorsqu'on modifie le commentaire.

### IV.1.2 Documentation

Maintenir une documentation à part sur le fonctionnement interne d'un programme est une mission quasiment impossible. C'est une méthode à éviter. Il vaut mille fois mieux intégrer la documentation au programme sous forme de commentaires.

Cette logique peut être poussée un peu plus loin en utilisant dans les commentaires les commandes d'un logiciel de formatage de documents (troff, L<sup>A</sup>T<sub>E</sub>X, etc.). Il suffit alors d'avoir un outil qui extrait les commentaires du source et les formate pour retrouver un document externe avec une présentation plus riche que des commentaires traditionnels. Knuth a formalisé cette approche sous le nom de programmation littéraire [12]

## IV.2 Typologie des noms

La typologie des noms (choix des noms des variables, des fonctions, des fichiers) est un élément primordial dans la lisibilité d'un programme. Celle-ci doit respecter plusieurs contraintes:

- **cohérence** choisissez une logique dans le choix des noms de variables et gardez-là.
- **signification** choisissez des noms qui ont un sens en relation avec le rôle de la variable ou de la fonction que vous nommez.
- **modularité** indiquez l'appartenance d'un nom à un module.
- **non-ambiguïté** évitez ambiguïtés pour distinguer deux variables semblables (variations sur la casse par exemple), utilisez des moyens simples (suffixes

numériques). Attention, certains éditeurs de liens imposent que les 6 (six!) premiers caractères d'un identificateur soient discriminants.

Il existe plusieurs conventions de choix des noms de variables et de fonctions décrites dans la littérature. Parmi celles, ci on peut citer la « notation hongroise » présentée entre autres par C. Simonyi [13] qui code le type des objets dans leur nom.

Sans entrer dans un mécanisme aussi systématique, il est bon de suivre quelques règles :

- les noms avec un underscore en tête ou en queue sont réservés au système. Ils ne doivent donc pas être utilisés par un utilisateur de base.
- mettre en majuscules les constantes et les noms de macros définies par `#define`. Les macros qui se comportent comme une fonction peuvent avoir un nom en minuscules.

exemples :

```
#define VITESSE_MAX (1.8)
#define MAX(i,j) ((i) > (j) ? (i) : (j))
#define bcopy(src,dst,n) memcpy((dst),(src),(n))
```

`MAX` est identifié comme une macro (et doit le rester). En effet cette macro évalue deux fois l'un de ses arguments. Écrire `max` risquerait de le faire oublier et de conduire à des erreurs.

- mettre en majuscules également les noms de types définis par `typedef` et les noms de structures.
- utiliser de préférence le même nom pour une structure et le type définit pour elle. exemple :

```
typedef struct POS {
    double x;
    double y;
    double theta;
} POS;
```

- les constantes dans les enum commencent par une majuscule.
- les autres noms (variables, fonctions,... ) commencent par une minuscule et sont essentiellement en minuscules. Quand un nom comporte plusieurs mots, on peut utiliser une majuscule pour introduire chaque nouveau mot.
- éviter les noms trop proches typographiquement. Par exemple les caractères «l» et «1» sont difficiles à distinguer, il en sera de même des identificateurs «u1» et «ul».

- si une fonction retourne une valeur qui doit être interprétée comme valeur booléenne dans un test, utiliser un nom significatif du test. Par exemple `valeurCorrecte()` plutôt que `testValeur()`.
- la longueur d'un nom n'est pas une vertu en soi. Un index de tableau n'a pas besoin d'être plus complexe que `i`. Les variables locales d'une fonction peuvent souvent avoir des noms très courts.
- les variables globales et les fonctions doivent au contraire avoir des noms qui comportent le maximum d'information. Mais attention, des noms trop longs rendent la lecture difficile.

## IV.3 Déclarations

Utilisez le C ANSI, et incluez systématiquement des prototypes des fonctions que vous utilisez. Tous les compilateurs C ANSI ont une option pour produire un avertissement ou une erreur quand une fonction est appelée sans que son prototype n'ait été déclaré.

Bien entendu, déclarez un type à toutes vos variables et à toutes vos fonctions. La déclaration implicite en entier est une source d'erreurs.

Pour déclarer des types compliqués, utilisez des `typedefs`. Cela rend le code plus lisible et plus modulaire.

## IV.4 Indentation

L'indentation permet de mettre en valeur la structure de l'algorithme. Il est capital de respecter une indentation cohérente avec cette structure. Mais, comme pour la typologie des noms de variables, il n'y a pas de règles uniques.

Personnellement, j'utilise un système d'indentation bien résumé par l'exemple suivant. Il a l'avantage d'une certaine compacité.

```
if (condition) {
    /* 1er cas */
    x = 2;
} else {
    /* 2nd cas */
    x = 3;
}
```

D'autres préfèrent aligner les accolades ouvrantes et fermantes qui se correspondent sur une même colonne :

```
if (condition)
```

```
{
    /* 1er cas */
    x = 2;
}
else
{
    /* 2nd cas */
    x = 3;
}
```

L'incrément de base de l'indentation doit être suffisant pour permettre de distinguer facilement les éléments au même niveau. Quatre caractères semble une bonne valeur.

Il existe plusieurs outils qui maintiennent l'indentation d'un programme automatiquement. L'éditeur `emacs` propose un mode spécifique pour le langage C qui indente les lignes tout seul au fur et à mesure de la frappe, selon des règles programmables.

L'utilitaire Unix `indent` permet de refaire l'indentation de tout un fichier. Un fichier de configuration permet de décrire son style d'indentation favori.

## IV.5 Boucles

Évitez absolument de transformer votre code en plat de spaghetti. Le langage C permet de nombreuses constructions qui détournent le cours normal de l'exécution du programme : `break`, `continue`, `goto`...

Toutes ces constructions doivent être évitées dès qu'elles rendent difficile le suivi du déroulement d'un programme. Par la suite, s'il faut prendre un compte un nouveau cas, cela ne pourra se faire qu'en ajoutant des nœuds dans le plat...

Mais attention, dans un certain nombre de cas, notamment le traitement des erreurs, l'utilisation judicieuse d'un `break` ou d'un `goto` est plus lisible qu'une imbrication profonde de tests.

## IV.6 Expressions complexes

Dé-com-po-sez les expressions trop complexes en utilisant éventuellement des variables intermédiaires. Cela diminue le risque d'erreur lors de la saisie et augmente la lisibilité pour la suite.

Pour déclarer un type complexe, utilisez plusieurs `typedefs` intermédiaires.

Par exemple, pour déclarer un tableau de dix pointeurs sur fonctions entières avec un paramètre entier, les deux typedefs suivants sont bien plus lisibles que ce que l'on obtiendrait en essayant de l'écrire directement (laissé en exercice pour le lecteur).

```
typedef int (*INTFUNC)(int);
typedef INTFUNC TABFUNC[10];
```

## IV.7 Conversion de types

Attention, terrain glissant. Normalement, il ne devrait pas y en avoir. Avant d'utiliser un cast, demandez-vous toujours s'il n'y a pas un problème dans votre programme qui vous oblige à faire ce cast.

Les pommes ne sont pas des poires, c'est vrai aussi des types informatiques. Si vraiment vous avez des types qui peuvent représenter plusieurs objets différents, les unions sont peut-être un peu plus lourdes à manier, mais elles offrent des possibilités de vérification au compilateur.

En effet, le plus grand piège tendu par les cast, est que vous obligez le compilateur à accepter ce que vous tapez, en lui ôtant tout droit à la critique. Or il est possible de faire des erreurs partout, y compris dans l'utilisation des cast. Mais le compilateur n'a plus aucun moyen de les détecter.

## IV.8 Assertions

Le mécanisme des assertions permet de déclarer des prédicats sur les variables d'une fonction qui doivent être vrais (appelés aussi *invariants*). En cas de situation anormale (en général à la suite d'une erreur de logique du programme) l'assertion fausse provoquera un arrêt du programme.

Les assertions sont introduites par le fichier d'en-tête `assert.h` et sont écrites sous la forme :

```
assert()(expression)
```

Ce mécanisme simple permet d'aider à la mise au point d'algorithmes un peu complexes, à la fois parce qu'ils guident le programmeur pendant le codage et qu'ils permettent d'aider à détecter les erreurs.

## Références bibliographiques

- [1] S. Summit. *C Programming FAQs: Frequently Asked Questions*. Addison-Wesley, 1995.
- [2] D. Goldberg. What every computer scientist should know about floating-point arithmetic. *ACM Computing Surveys*, 23(1):5–48, March 1991.
- [3] D.E. Knuth. *Seminumerical Algorithms*, volume 2 of *The Art of Computer Programming*. Addison-Wesley, 1973.
- [4] B.W. Kernighan and D.M. Ritchie. *The C Programming Language*. Prentice-Hall, 1978.
- [5] B.W. Kernighan and D.M. Ritchie. *The C Programming Language*. Prentice-Hall, 2nd edition, 1988.
- [6] D.E. Knuth. *Fundamental Algorithms*, volume 1 of *The Art of Computer Programming*. Addison-Wesley, 1973.
- [7] R. Pike. *Notes on Programming in C*.
- [8] D.E. Knuth. *Sorting and Searching*, volume 3 of *The Art of Computer Programming*. Addison-Wesley, 1973.
- [9] Aleph One ([aleph1@underground.org](mailto:aleph1@underground.org)). Smashing the stack for fun and profit. *Phrack*, (49), November 1996.
- [10] S. Garfinkel and G. Spafford. *Practical Unix and Internet Security*. O'Reilly and Associates, 2nd edition, 1996.
- [11] L.W. Cannon, R.A. Elliot, L.W. Kirchhoff, J.H. Miller, J.M. Milner, R.W. Mitze, E.P. Shan, and N.O. Whittington. *Indian Hill C style and coding standards*. Bell Labs.
- [12] D.E. Knuth. Literate programming. *Computer Journal*, 28(2):97–111, 1984.
- [13] C. Simonyi and M. Heller. The hungarian revolution. *Byte*, pages 131–138, Août 1991.

# Index

## Symboles

<code>&amp;</code> .....	25
<code>&amp;&amp;</code> .....	13
<code>=</code> .....	7
<code>==</code> .....	7
<code> </code> .....	25
<code>  </code> .....	13
<code>~</code> .....	25
<code>-</code> .....	20
64 bits .....	19

## A

allocation mémoire .....	13
ANSI .....	11
arrondi .....	10
<code>assert</code> .....	41
<code>assert.h</code> .....	41
assertions .....	41

## B

boucles .....	40
<code>break</code> .....	8, 40
<code>bsearch</code> .....	25

## C

calcul réel .....	9
caractères	
chaînes de .....	15
<code>case</code> .....	8
<code>ceil</code> .....	10
<code>char</code> .....	12, 15
commentaires .....	35
<code>const</code> .....	16

<code>continue</code> .....	40
<code>cpp</code> .....	20

## D

déclarations .....	11, 39
documentation .....	37
données	
binaires .....	19
<code>double</code> .....	10, 11, 20

## E

égalité .....	9
<code>emacs</code> .....	40
en-tête .....	35
entrées/sorties .....	17
<code>exec</code> .....	33

## F

FAQ .....	5
<code>fgets</code> .....	15, 18, 19, 33
<code>float</code> .....	12
<code>floor</code> .....	10
<code>for</code> .....	14
<code>fork</code> .....	33
<code>fprintf</code> .....	18
<code>fread</code> .....	19
<code>free</code> .....	13, 14, 17
<code>fscanf</code> .....	18
fuites .....	15
<code>fwrite</code> .....	19

## G

<code>getchar</code> .....	18, 19
----------------------------	--------

- gets ..... 15, 19, 33  
goto ..... 40
- H**  
hsearch ..... 25
- I**  
indent ..... 40  
indentation ..... 39  
int ..... 10, 20
- K**  
Kernigan et Ritchie ..... 11
- L**  
lex ..... 26  
long int ..... 11  
lsearch ..... 25
- M**  
malloc ..... 13, 15, 17, 23, 30  
math.h ..... 10  
mkstemp ..... 33  
mktemp ..... 33
- N**  
non-initialisées  
    variables ..... 12
- O**  
ordre d'évaluation ..... 12
- P**  
passage par adresse ..... 9  
pointeurs ..... 17  
préprocesseur ..... 20  
printf ..... 18
- Q**  
qsort ..... 25
- S**  
sbrk ..... 23  
scanf ..... 9, 18, 19, 33  
short ..... 12  
snprintf ..... 16, 33  
sprintf ..... 15, 16, 18, 33  
sscanf ..... 18  
stdlib.h ..... 10  
strcat ..... 33  
strcpy ..... 15, 16, 33  
strncat ..... 33  
strncpy ..... 16, 33  
strtod ..... 10, 33  
strtok ..... 33  
strtol ..... 33  
switch ..... 8  
system ..... 31, 33
- T**  
tableaux ..... 7, 17  
typedef ..... 38, 39  
types  
    conversion de ..... 41  
typologie ..... 37
- U**  
union ..... 20  
unsigned int ..... 20  
Usenet ..... 5
- Y**  
yacc ..... 26