

ALGORITHME LZW

1 Structures de données

Conceptuellement, on a besoin de deux structures de données pour implémenter l'algorithme LZW :

- une structure (qu'on appellera dictionnaire) donnant les associations code vers motif ;
- une structure (qu'on appellera table inverse) donnant les associations motif vers code (et permettant de tester si un motif est associé à un code).

On peut imaginer la solution suivante :

- la première structure (le dictionnaire) est réalisée par un tableau, avec le motif associé au code c dans la case d'indice c ;
- la seconde structure (la table inverse) est réalisée par une table de hachage, avec les motifs comme clés et les codes comme valeurs.

► **Question 1** En supposant qu'on choisisse cette solution :

- les deux structures sont-elles utiles dans la phase de compression ?
- et dans la phase de décompression ?

La solution que nous allons utiliser est légèrement différente, et tire parti du fait que la structure des motifs « connus » (*i.e* associés à un code) est très contrainte. Un motif présent dans la table est :

- soit réduit à un octet ;
- soit de la forme mx , avec m un autre motif présent dans le dictionnaire et x un octet.

On peut donc éviter de manipuler des motifs en tant que tels, et n'utiliser que des couples c, x où :

- c est un code déjà existant ;
- x est un octet.

On définit les alias de type suivants :

```
// codeword type
typedef uint32_t cw_t;
// byte type
typedef uint8_t byte_t;
```

Une entrée du dictionnaire est un couple (code, octet) :

```
struct dict_entry_t {
    cw_t pointer;
    byte_t byte;
};

typedef struct dict_entry_t dict_entry_t;
```

Le dictionnaire en lui-même est une structure contenant :

- le prochain code disponible ;
- la largeur d'un code (en nombre de bits) – pour l'instant, cette taille sera constante ;
- un tableau statique de `dict_entry_t`. La taille de ce tableau est une constante globale, et vaut 2^d , où d est la largeur maximale d'un code.

```

#define CW_MAX_WIDTH 16
#define DICTFULL (1u << CW_MAX_WIDTH)

const cw_t NO_ENTRY = DICTFULL;
const cw_t NULL_CW = DICTFULL;
const cw_t FIRST_CW = 0x100;
const int CW_MIN_WIDTH = 16;

struct dict_entry_t {
    cw_t pointer;
    byte_t byte;
};

typedef struct dict_entry_t dict_entry_t;

struct dict_t {
    cw_t next_available_cw;
    int cw_width;
    dict_entry_t data[DICTFULL];
};

struct dict_t dict;

```

Remarques

- On utilise la directive du pré-processeur `#define` pour les constantes `CW_MAX_WIDTH` et `DICTFULL` pour pouvoir définir un tableau statique de taille `DICTFULL` (ce ne serait pas possible si `DICTFULL` était défini comme un `const int`, par exemple). Ce point peut être ignoré.
- `dict` est une variable globale.
- `NO_ENTRY` et `NULL_CW` sont des constantes dont on sait qu'elles ne peuvent correspondre à un code valide. On a trois constantes différentes pour la même valeur 2^d , mais elles seront utilisées dans des contextes différents.
- `FIRST_CW` indique le premier code créé dynamiquement (après les codes pour les motifs de un octet). On a donné sa valeur en hexadécimal, qui correspond à 256 en décimal.

► **Question 2** Avec les valeurs données ci-dessus pour les différentes constantes (et pour le type `cw_t`), quelle quantité de mémoire le dictionnaire occupe-t-il ?

► **Question 3** Écrire une fonction `initialize_dictionary` qui initialise les champs `next_available_cw` et `cw_width`, ainsi que la partie du tableau correspondant aux motifs de un octet.

- Pour `cw_width`, on initialisera à `CW_MIN_WIDTH`.
- On mettra le champ `pointer` à `NULL_CW` pour les motifs de un octet.

```
void initialize_dictionary(void);
```

Pour la table inverse (association motif vers code), on utilise un tableau bidimensionnel de codes :

```
cw_t inverse_table[DICTFULL][256];
```

La case `inverse_table[c][x]` contiendra le code correspondant au couple (c, x) s'il existe, une valeur quelconque sinon.

► **Question 4** Quelle quantité de mémoire `inverse_table` consomme-t-elle ?

► **Question 5** Écrire une fonction `lookup` qui prend en entrée un couple `c, x` et renvoie :

- le code correspondant à `c, x` s'il y en a un ;
- `NO_ENTRY` sinon.

```
cw_t lookup(cw_t cw, byte_t byte);
```

► **Question 6** Écrire une fonction `build_entry` qui prend en entrée un couple `c, x` et ajoute une entrée dans le dictionnaire pour ce motif (en mettant également à jour `inverse_table`). On pourra supposer sans le vérifier que le motif n'est pas déjà présent dans le dictionnaire.

Remarque

Dans le cas où le dictionnaire est déjà plein, cette fonction ne fera rien.

```
void build_entry(cw_t cw, byte_t byte);
```

2 Compression

Pour compresser, nous allons lire le fichier d'entrée octet par octet et émettre des codes sur un fichier de sortie. Pour l'instant, ces codes seront émis en ASCII : si l'on émet 517, on écrira "517" (trois caractères ASCII) sur le fichier de sortie. Bien sûr, cela ne permet pas de compresser réellement, mais cela nous permettra de vérifier la correction de notre algorithme de compression.

► **Question 7** Écrire une fonction `mock_compress` qui prend en entrée un fichier d'entrée et un fichier de sortie et écrit dans le fichier de sortie les codes générés, au format ASCII, séparées par une espace. Par exemple, avec en entrée un fichier réduit à "ABABCABBAB\n", et sachant que le code ASCII de A est 65, on doit obtenir le résultat suivant : "65 66 256 67 256 257 66 10".

```
void mock_compress(FILE *input_file, FILE *output_file);
```

Remarque

Pour lire un octet du fichier d'entrée, on utilisera la fonction `getc`. Son prototype est :

```
int getc(FILE *stream);
```

L'entier qu'elle renvoie est soit EOF (une constante prédéfinie, strictement négative, qui signifie qu'on est arrivé à la fin du fichier), soit une valeur entre 0 et 255 (qui peut donc être transtypée sans problème en `byte_t`).

► **Question 8** Écrire un programme ayant le comportement suivant :

- il accepte entre un et trois arguments en ligne de commande ;
- le premier argument est réduit à un caractère :
 - `c` pour compresser en mode binaire à l'aide de la fonction `compress` (qui reste à écrire) ;
 - `C` pour compresser en mode texte avec `mock_compress` ;
 - `d` pour décompresser en mode binaire ;
 - `D` pour décompresser en mode texte ;
- le deuxième argument, s'il est présent, indique le fichier d'entrée (sinon, on prend l'entrée standard) ;
- le troisième argument, s'il est présent, indique le fichier de sortie (sinon, on prend la sortie standard).

Remarque

On en profitera bien entendu pour tester la fonction `mock_compress` sur l'exemple ci-dessus, et sur d'autres exemples de préférence !

► **Question 9** Quelle est la complexité temporelle totale de votre programme, en fonction de la largeur `d` des codes et du nombre `n` d'octets à compresser ?

3 Décompression

La décompression est un peu plus délicate que la compression à cause de la manière dont nous avons choisi de représenter le dictionnaire. Implicitement, le dictionnaire contient les motifs sous forme de listes chaînées de caractères, avec le dernier caractère du motif en tête de liste. Nous voulons bien sûr émettre les caractères dans l'ordre, ce qui peut se faire de deux manières :

- soit à l'aide d'une pile;
- soit à l'aide d'une fonction récursive.

La version utilisant une fonction récursive est légèrement plus facile à écrire, mais peut être problématique dans les cas pathologiques puisque elle utilise un espace proportionnel à la longueur du motif sur la pile d'appel. On fournit donc une réalisation très simple de la structure de pile *via* le header `stack.h`, qui déclare les fonctions suivantes :

```
typedef struct stack stack;

stack *stack_new(int capacity);
void stack_free(stack *s);
int stack_size(stack *s);
byte_t stack_pop(stack *s);
void stack_push(stack *s, byte_t byte);
```

► **Question 10** Écrire une fonction `decode_cw` ayant le prototype suivant :

```
byte_t decode_cw(FILE *fp, cw_t cw, stack *s);
```

Cette fonction émettra, dans l'ordre, tous les octets du motif dont le code est `cw` sur le fichier `*fp`. La pile est fournie pour éviter de la réallouer à chaque appel : on pourra supposer qu'elle est de taille suffisante pour contenir tous les octets du motif, et son état tant avant qu'après l'appel n'a pas d'importance. De plus, elle renverra le dernier octet du motif (qui nous sera utile par la suite). Pour émettre un octet sur le flux de sortie, on utilisera :

```
int putc(int ch, FILE *stream);
```

L'argument `ch` est automatiquement converti en `unsigned char` avant d'être écrit, on donnera donc directement un `byte_t` comme argument. La valeur de retour est égale à `ch` si tout s'est bien passé, à `E0F` sinon : on pourra l'ignorer.

► **Question 11** Écrire une fonction `get_first_byte` qui renvoie le premier octet du motif associé à un code.

```
byte_t get_first_byte(cw_t cw);
```

► **Question 12** Écrire une fonction `mock_decompress` qui lit un fichier compressé au format produit par `mock_compress` et écrit le flux décompressé correspondant sur `output_file`.

- Il est conseillé de commencer par retrouver l'algorithme de décompression par soi-même : il n'est pas invisable que l'on vous demande cela le jour d'un concours...
- Après avoir fourni cet effort, consulter le cours est quand même une bonne idée.
- La fonction `decode_cw` suffit pour traiter le cas « usuel » ; la fonction `get_first_byte` est utile pour traiter le cas dit KwK (celui où l'on lit un code que l'on n'a pas encore ajouté au dictionnaire).
- La table inverse ne sert pas pour la décompression.

```
void mock_decompress(FILE *input_file, FILE *output_file);
```

On pensera à tester la fonction sur une entrée contenant un cas KwK (par exemple "ABABABA\n")!

4 Lecture et écriture binaires

Pour réaliser une véritable compression, il est nécessaire qu'un code de largeur d utilise d bits sur le fichier de sortie. Comme d n'a aucune raison d'être un multiple de 8, on est ramené à un problème similaire à celui que l'on a résolu en OCaml pour le code de Huffman. Nous allons procéder de manière très légèrement différente :

- on maintiendra toujours un accumulateur (que nous appellerons `buffer`) et sa taille (en nombre de bits significatifs), et l'on écrira toujours un octet sur le fichier de sortie quand le nombre de bits de l'accumulateur atteindra ou dépassera 8;
- cependant, au lieu de recevoir les bits à écrire un par un, nous les recevrons par paquet (un code complet à chaque appel);
- d'autre part, la clôture du fichier sera plus simple : on pourra se contenter de compléter le dernier octet par des zéros. En effet, lors de la décompression, on saura toujours combien de bits on souhaite lire, et ce nombre sera toujours supérieur ou égal à 9 (longueur minimale possible d'un code).

On utilise la structure suivante :

```
const int BUFFER_WIDTH = 64;
const int BYTE_WIDTH = 8;
const uint64_t BYTE_MASK = (1 << BYTE_WIDTH) - 1;

struct bit_file {
    FILE *fp;
    uint64_t buffer;
    int buffer_length;
};

typedef struct bit_file bit_file;

bit_file *bin_initialize(FILE *fp){
    bit_file *bf = malloc(sizeof(bit_file));
    bf->fp = fp;
    bf->buffer = 0;
    bf->buffer_length = 0;
    return bf;
}
```

► **Question 13** En utilisant un *buffer* de 64 bits, quelle taille de code peut-on traiter au maximum sans problème? La limitation est-elle gênante?

► **Question 14** Écrire une fonction `output_bits` de prototype :

```
void output_bits(bit_file *bf, uint64_t data, int width, bool flush);
```

Les données à écrire sont des `width` bits de poids faible de `data`. Le paramètre `flush` détermine le comportement sur les bits restants après avoir écrit autant d'octets que possible dans `bf` :

- s'il vaut `false`, les bits restants sont laissés dans `bf->buffer`;
- s'il vaut `true`, ils sont écrits dans `bf`, complétés par des zéros pour obtenir un octet.

On écrira les données bit le moins significatif en premier : en particulier, quand on écrit un octet constitué du reste du code précédent et du début du code actuel, les bits de poids faible de l'octet correspondront à l'ancien code et ceux de poids fort au nouveau.

► **Question 15** Écrire une fonction `input_bits` ayant le prototype suivant :

```
uint64_t input_bits(bit_file *bf, int width, bool *eof);
```

Cette fonction lit width bits depuis le flux bit_file (de manière à lire correctement un flux écrit par output_bits, évidemment). La valeur pointée par eof sera mise à `true` si la lecture a échoué parce que l'on est arrivé à la fin du fichier sans parvenir à lire width bits, à `false` sinon.

Remarque

Les width bits lus seront placés dans les bits de poids faibles de la valeur de retour.

► **Question 16** Écrire les fonctions `compress` et `decompress`, ayant les mêmes prototypes que `mock_compress` et `mock_decompress` mais écrivant les codes en version « compacte ».

```
void compress(FILE *input_file, FILE *output_file);
void decompress(FILE *input_file, FILE *output_file);
```

► **Question 17** Tester le taux de compression obtenu pour :

- le fichier source de votre code C d'aujourd'hui;
- l'énoncé du TP en format pdf;
- l'exécutable obtenu en compilant votre code source;
- le texte intégral de *Moby Dick* fourni avec le sujet.

On pourra comparer :

- les taux de compression obtenus pour différentes largeurs de code;
- le taux de compression obtenus en utilisant l'utilitaire zip.

5 Codes de largeur variable

Pour améliorer le taux de compression, une solution simple est d'utiliser des codes à largeur variable. En effet, en supposant que l'on fixe la largeur des codes à 14 bits par exemple, on va mettre assez longtemps à émettre le premier code ne rentrant pas sur 13 bits (*i.e.* 8 192) : jusque là, le ou les bits les plus significatifs des codes émis valent tous zéro, et l'on gaspille donc de la place.

Pour éviter cela, il suffit de se mettre d'accord (entre la fonction de compression et celle de décompression) sur une règle pour l'évolution de la largeur du code. La règle la plus simple est la suivante :

- au départ, la largeur d'un code est 9 bits;
- on se fixe une largeur maximale (disons 16 bits), et l'on crée les structures `dict` et `inverse_table` avec une taille correspondant à cette largeur;
- dès que l'on souhaite créer une entrée pour un code et que ce code ne tient pas sur le nombre actuel de bits, on augmente la largeur de 1 si c'est possible (sinon, le dictionnaire est plein et l'entrée n'est pas créée);
- le nouveau code est créé au moment où l'on émet un code déjà existant (systématiquement) : on convient que le code existant est émis avec l'ancienne largeur, ce qui revient à dire que l'on crée l'entrée pour le nouveau code après émission de l'ancien;
- pour la décompression, il faut juste penser que l'on a toujours « un temps de retard » sur la compression (pour l'état du dictionnaire), et qu'il faut donc changer de largeur une étape plus tôt.

► **Question 18** Modifier la fonction `build_entry` pour qu'elle mette à jour la largeur des codes. Le comportement étant légèrement différent suivant que l'on est en train de compresser ou de décompresser, on ajoute un paramètre booléen `compress_mode` pour indiquer le mode de fonctionnement.

```
void build_entry(cw_t cw, byte_t byte, bool compress_mode);
```

► **Question 19** Après avoir apporté les autres modifications nécessaires à votre code (s'il y a lieu), reprendre les mesures de taux de compression et les comparaisons avec zip. On pourra aussi comparer avec la compression de Huffman que nous avons déjà programmé, et tester si la composée de Huffman et de LZW, dans un sens ou dans l'autre, présente un intérêt.

Solutions

► **Question 1** Pour la phase de compression, on n'a besoin que de l'association motif vers code, donc du dictionnaire. Pour la phase de décompression, c'est l'inverse : il suffit d'avoir l'association code vers motif.

► **Question 2**

Réponse attendue : une `dict_entry_t` occupe $4 + 1 = 5$ octets, donc le dictionnaire occupe $5 \cdot 2^{16} \simeq 320\text{KiB}$ (en négligeant les quelques octets occupés par les champs `next_available_cw` et `cw_width`).

En réalité : pour des raisons d'alignement, un tableau de n `dict_entry_t` occupe en fait $8n$ octets. On obtient donc $2^{19} = 512\text{KiB}$.

► **Question 3** Aucune difficulté, juste une remarque : les constantes globales sont là pour être utilisées. En général, on considère que les seules constantes littérales qui peuvent apparaître directement dans le code (sans qu'on leur ait donné un nom) sont $-1, 0, 1$ et 2 .

```
void initialize_dictionary(void){
    dict.next_available_cw = FIRST_CW;
    dict.cw_width = CW_MIN_WIDTH;
    for (cw_t i = 0; i < FIRST_CW; i++) {
        dict.data[i].pointer = NULL_CW;
        dict.data[i].byte = i;
    }
}
```

► **Question 4** `inverse_table` est un tableau de $2^{16} \cdot 2^8 = 2^{24}$ `cw_t`, qui occupent chacun 4 octets. Au total, on a donc $2^{28} = 256\text{MiB}$. Cela commence à ne pas être négligeable : en limitant la largeur des codes à 15 bits, et en utilisant `uint16_t` pour le type `cw_t` (ce qui serait alors possible), on diviserait la consommation mémoire par 4 et n'utiliserait plus que 64MiB .

► **Question 5** `inverse_table` indique l'unique case de `dict` dans laquelle l'entrée cherchée peut se trouver. On vérifie ensuite le contenu de cette case, puisqu'on nous dit qu'on ne peut rien supposer sur le contenu de `inverse_table` pour les codes absents.

```
cw_t lookup(cw_t cw, uint8_t byte){
    cw_t address = inverse_table[cw][byte];
    if (address >= dict.next_available_cw) return NO_ENTRY;
    dict_entry_t entry = dict.data[address];
    if (entry.pointer == cw && entry.byte == byte) return address;
    return NO_ENTRY;
}
```

On observe que cette fonction s'exécute en $O(1)$.

► **Question 6** Pas de difficulté particulière :

```

void build_entry(cw_t cw, uint8_t byte){
    cw_t next = dict.next_available_cw;

    if (next == DICTFULL) return;

    dict.data[next].pointer = cw;
    dict.data[next].byte = byte;

    inverse_table[cw][byte] = next;

    dict.next_available_cw++;
}

```

► Question 7

```

1 void mock_compress(FILE *input_file, FILE *output_file){
2     cw_t current_cw = NULL_CW;
3     int current_byte = getc(input_file);
4     while (current_byte != EOF) {
5         if (current_cw == NULL_CW) current_cw = current_byte;
6         else {
7             cw_t new_cw = lookup(current_cw, (uint8_t)current_byte);
8             if (new_cw == NO_ENTRY) {
9                 fprintf(output_file, "%d ", current_cw);
10                build_entry(current_cw, (uint8_t)current_byte, true);
11                current_cw = (cw_t)current_byte;
12            } else {
13                current_cw = new_cw;
14            }
15        }
16        current_byte = getc(input_file);
17    }
18    fprintf(output_file, "%d", current_cw);
19    if (VERBOSITY > 0) {
20        fprintf(stderr, "Distinct codewords : %d\n", (int)dict.next_available_cw);
21    }
22 }

```

Quelques remarques :

- le test ligne 5 sert juste à traiter le premier caractère;
- le `fprintf` après la boucle sert à émettre le code du motif actuel lorsque l'on atteint la fin du fichier;
- on affiche (suivant la valeur de la variable globale `VERBOSITY`) le nombre de codes créés sur la sortie d'erreur standard.

► Question 8 Essayons de faire les choses à peu près proprement :

- on définit une fonction `print_usage_and_exit` qui affiche un message d'aide et quitte le programme (grâce à un appel à `exit(EXIT_FAILURE)`);
- dans la fonction `main`, on vérifie que la ligne de commande correspond à ce que l'on attend (entre un et trois arguments, avec le premier argument réduit à un caractère);
- si c'est le cas, on tente d'ouvrir les fichiers spécifiés, en signalant une erreur éventuelle;
- finalement, on effectue l'action demandée puis l'on ferme les fichiers (fermer `stdin` ou `stdout` est légal tant qu'on ne s'en sert pas par la suite).


```

void print_usage_and_exit(char *command){
    fprintf(stderr, "Usage :\n");
    fprintf(stderr, "  %s c <input-file> <output-file>"
        " to compress in binary mode\n", command);
    fprintf(stderr, "  %s C <input-file> <output-file>"
        " to compress in ascii mode\n", command);
    fprintf(stderr, "  %s d <input-file> <output-file>"
        " to decompress in binary mode\n", command);
    fprintf(stderr, "  %s D <input-file> <output-file>"
        " to decompress in ascii mode\n", command);
    fprintf(stderr, "If only one file is given, output is to stdout.\n");
    fprintf(stderr, "If no file is given, input is from stdin"
        " and output is to stdout.\n");
    exit(EXIT_FAILURE);
}

int main(int argc, char* argv[]){
    FILE *input_file = stdin;
    FILE *output_file = stdout;
    if (argc < 2 || argc > 4 || strlen(argv[1]) != 1) {
        print_usage_and_exit(argv[0]);
    }
    if (argc >= 3) input_file = fopen(argv[2], "rb");
    if (argc >= 4) output_file = fopen(argv[3], "wb");
    if (input_file == NULL) {
        fprintf(stderr, "Cannot open file %s for reading.\n", argv[2]);
        return(EXIT_FAILURE);
    }
    if (output_file == NULL) {
        fprintf(stderr, "Cannot open file %s for writing.\n", argv[2]);
        return(EXIT_FAILURE);
    }

    initialize_dictionary();
    char action = argv[1][0];
    if (action == 'c') compress(input_file, output_file);
    else if (action == 'C') mock_compress(input_file, output_file);
    else if (action == 'd') decompress(input_file, output_file);
    else if (action == 'D') mock_decompress(input_file, output_file);
    else print_usage_and_exit(argv[0]);

    fclose(input_file);
    fclose(output_file);
    return EXIT_SUCCESS;
}

```

Remarques

- Attention, on peut comparer `argv[1][0]` à `'c'` (par exemple) avec `==`, mais pas `argv[1]` à `"c"`. Dans le deuxième cas, on comparerait des `char*` et l'on testerait donc l'égalité des *pointeurs* ! Si l'on veut comparer deux chaînes, il faut utiliser la fonction `strcmp` (fournie dans `string.h`).
- Un `switch` serait clairement plus idiomatique pour la disjonction de cas sur `action`, mais ce n'est pas au programme et la version proposée n'est pas tellement plus lourde.

► **Question 9** On va considérer que la création de `dict` et `inverse_table` se font en temps proportionnel à leur taille, donc $O(2^d)$. Ensuite, on traite les octets un par un avec à chaque fois un appel à `lookup` et éventuellement un appel à `build_entry` : ces deux fonctions étant en temps constant, on obtient $O(n)$ pour la phase de compression proprement dite. Au total, la complexité est donc de $O(n + 2^d)$.

► **Question I0** On empile tous les octets du code puis on les émet en les dépilant. On retient le dernier octet empilé (et donc le premier émis) pour le renvoyer en valeur de retour.

```
uint8_t decode_cw(FILE *fp, cw_t cw, stack *s){
    assert(cw < dict.next_available_cw);
    dict_entry_t entry = dict.data[cw];
    while (entry.pointer != NULL_CW) {
        stack_push(s, entry.byte);
        entry = dict.data[entry.pointer];
    }
    stack_push(s, entry.byte);

    while (stack_size(s) > 0) {
        putc(stack_pop(s), fp);
    }

    return entry.byte;
}
```

► **Question I1** Il s'agit essentiellement de trouver le dernier élément d'une liste simplement chaînée :

```
uint8_t get_first_byte(cw_t cw){
    dict_entry_t entry = dict.data[cw];
    while (entry.pointer != NULL_CW) {
        entry = dict.data[entry.pointer];
    }
    return entry.byte;
}
```

► **Question I2**

```
void mock_decompress(FILE *input_file, FILE *output_file){
    stack *s = stack_new(DICTFULL);
    cw_t prev_cw;
    cw_t current_cw;
    if (fscanf(input_file, "%d", &prev_cw) == 1) {
        decode_cw(output_file, prev_cw, s);
    }
    while (fscanf(input_file, "%d", &current_cw) == 1) {
        if (current_cw < dict.next_available_cw) {
            uint8_t byte = decode_cw(output_file, current_cw, s);
            build_entry(prev_cw, byte, false);
        } else {
            // KwK case
            uint8_t byte = get_first_byte(prev_cw);
            build_entry(prev_cw, byte, false);
            decode_cw(output_file, current_cw, s);
        }
        prev_cw = current_cw;
    }
    stack_free(s);
}
```

Remarques

- On rappelle que `fscanf` renvoie le nombre de valeurs lues avec succès, ce qui permet de détecter la fin du fichier.

- Le premier code doit être traité à part (puisque `prev_cw` n'est pas encore défini).
- On crée une pile de taille `DICTFULL`, ce qui est clairement suffisant pour n'importe quel code. Dans le cas d'un texte constitué d'un seul octet répété `n` fois, on peut (presque) atteindre cette borne.
- On n'oublie bien sûr pas de libérer la pile avant de sortir de la fonction.

► **Question I3** Le *buffer* contient au plus 7 bits au début d'un appel (puisque on écrit un octet dès que l'on dispose de 8 bits). Avec 64 bits, on peut donc accepter sans problème des codes jusqu'à 57 bits : c'est très largement suffisant, puisqu'un dictionnaire de taille 2^{57} est inenvisageable. En réalité, un *buffer* de 32 bits, permettant des codes de 25 bits, serait largement suffisant.

► **Question I4**

```
void output_bits(bit_file *bf, uint64_t data, int width, bool flush){
    assert(bf->buffer_length + width <= BUFFER_WIDTH);
    data &= (1 << width) - 1;
    bf->buffer |= (data << bf->buffer_length);
    bf->buffer_length += width;
    while (bf->buffer_length >= BYTE_WIDTH) {
        fputc(bf->buffer & BYTE_MASK, bf->fp);
        bf->buffer >>= BYTE_WIDTH;
        bf->buffer_length -= BYTE_WIDTH;
    }
    if (flush && bf->buffer_length > 0) {
        fputc(bf->buffer & BYTE_MASK, bf->fp);
        bf->buffer = 0;
        bf->buffer_length = 0;
    }
}
```

Remarque

`&=`, `|=`, `>>=` ont un sens qui ne devrait pas être difficile à deviner. On peut bien sûr s'en passer.

► **Question I5**

```
uint64_t input_bits(bit_file *bf, int width, bool *eof){
    int byte = 0;
    int offset = bf->buffer_length;
    while (byte != EOF && bf->buffer_length < width) {
        byte = getc(bf->fp);
        bf->buffer |= (byte & BYTE_MASK) << offset;
        bf->buffer_length += BYTE_WIDTH;
        offset += BYTE_WIDTH;
    }
    if (byte == EOF) {
        *eof = true;
        return 0;
    }
    uint64_t buffer_mask = (1 << width) - 1;
    uint64_t res = bf->buffer & buffer_mask;
    bf->buffer >>= width;
    bf->buffer_length -= width;
    return res;
}
```

Remarque

Même si l'on a réussi à lire des bits, si l'on obtient EOF avant d'avoir réussi à lire le nombre demandé de bits, cela signifie qu'il n'y a plus de codes dans le fichier. Il faut donc jeter le résultat et signaler la fin de la lecture en mettant `*eof` à `true`.

► **Question I6** Les modifications à apporter à `mock_compress` sont triviales; Pour `decompress`, il y a un tout petit peu plus de travail puisque `input_bits` a un prototype assez différent de `fscanf`. Rien de bien compliqué toutefois.

```
void compress(FILE *input_file, FILE *output_file){
    bit_file *bf = bin_initialize(output_file);
    cw_t current_cw = NULL_CW;
    int current_byte = getc(input_file);
    while (current_byte != EOF) {
        if (current_cw == NULL_CW) current_cw = current_byte;
        else {
            cw_t new_cw = lookup(current_cw, (uint8_t)current_byte);
            if (new_cw == NO_ENTRY) {
                output_bits(bf, current_cw, dict.cw_width, false);
                build_entry(current_cw, (uint8_t)current_byte, true);
                current_cw = (cw_t)current_byte;
            } else {
                current_cw = new_cw;
            }
        }
        current_byte = getc(input_file);
    }
    output_bits(bf, current_cw, dict.cw_width, true);
    if (VERBOSITY > 0) {
        fprintf(stderr, "Distinct codewords : %d\n", (int)dict.next_available_cw);
    }
    free(bf);
}

void decompress(FILE *input_file, FILE *output_file){
    bit_file *bf = bin_initialize(input_file);
    stack *s = stack_new(DICTFULL);
    bool done = false;

    cw_t prev_cw = input_bits(bf, dict.cw_width, &done);
    if (!done) {
        decode_cw(output_file, prev_cw, s);
    }
    cw_t current_cw = input_bits(bf, dict.cw_width, &done);
    while (!done) {
        if (current_cw < dict.next_available_cw) {
            uint8_t byte = decode_cw(output_file, current_cw, s);
            build_entry(prev_cw, byte, false);
        } else {
            if (VERBOSITY >= 1) {
                fprintf(stderr, "KwK case, cw = %d\n", (int)current_cw);
            }
            uint8_t byte = get_first_byte(prev_cw);
            build_entry(prev_cw, byte, false);
            decode_cw(output_file, current_cw, s);
        }
        prev_cw = current_cw;
        current_cw = input_bits(bf, dict.cw_width, &done);
    }
    stack_free(s);
    free(bf);
}
```

► **Question 17** Se référer à la dernière question.

► **Question 18** Il faut juste faire attention au décalage de un entre compression et décompression.

```
void build_entry(cw_t cw, uint8_t byte, bool compress_mode){
    cw_t next = dict.next_available_cw;
    cw_t next_growth = 1u << dict.cw_width;

    if (next == DICTFULL) return;
    if (compress_mode && next == next_growth) dict.cw_width++;
    if (!compress_mode && next + 1 == next_growth) dict.cw_width++;

    dict.data[next].pointer = cw;
    dict.data[next].byte = byte;

    inverse_table[cw][byte] = next;

    dict.next_available_cw++;
}
```

► **Question 19** Il n'y a normalement rien d'autre à changer. On donne le taux de compression (taille finale / taille initiale) obtenu pour :

- `lzw_variable.c`, le code source du corrigé;
- `moby_dick.txt`, le texte intégral de *Moby Dick*;
- `variable`, l'exécutable obtenu en compilant `lzw_variable.c` (plus les deux fichiers dont il dépend).

	Taille initiale	LZW largeur fixe				LZW variable	Huffman	zip
		10	12	14	16			
<code>lzw_variable.c</code>	7.7ko	1.71	2.08	1.79	1.57	2.33	1.57	4.05
<code>variable</code>	117ko	1.95	1.86	2.66	2.60	3.00	2.09	4.68
<code>moby_dick.txt</code>	1.3Mo	1.53	1.82	2.15	2.34	2.37	1.69	2.43

Remarque

zip utilise l'algorithme deflate, qui combine LZ77 (une autre méthode de compression par dictionnaire) avec le codage de Huffman. On voit qu'il donne systématiquement de meilleurs résultats que LZW ou Huffman utilisés seuls, ce qui est dû à plusieurs facteurs :

- pour les petits fichiers, deflate peut utiliser un arbre de Huffman prédéfini (qui n'est donc pas adapté spécifiquement au fichier d'entrée, mais qui permet d'économiser l'espace de stockage de l'arbre);
- l'algorithme de Huffman est légèrement modifié dans deflate pour bien s'adapter au résultat de LZ77.

Combiner la version standard de Huffman avec LZW ne permet pas d'obtenir un meilleur taux de compression que LZW seul (que l'on commence par LZW – le plus logique – ou par Huffman).