

# Rapport de TP 4MMAOD : Génération de patch optimal

CASTRO LOPEZ Walter Ferney<sub>1</sub> (groupe 7<sub>1</sub>)  
BUALO GIORDANO Lucas Santiago<sub>2</sub> (groupe 7<sub>2</sub>)

28 novembre 2015

## 1 Principe de notre programme

Nous avons implémenté notre programme avec la méthode récursive. Pour cela, nous avons créé un tableau de dimension  $n \times m$  où  $n$  est le nombre de lignes du fichier original (d'entrée) et  $m$  est le nombre de lignes du fichier résultat (de sorti); et dans laquelle nous gardons les valeurs des coûts optimaux pour calculer les patches qui transforment les fichiers  $A_i$  en  $B_j$ , c'est-à-dire, le fichier qui va de la ligne 1 jusqu'à la ligne  $i$  du fichier d'entrée, et le fichier qui va de la ligne 1 jusqu'à la ligne  $j$  du fichier de sorti respectivement.

Autre que le coût, on mémorise pour chaque patch optimal la dernière instruction qui a été effectuée, permettant de retrouver le patch optimal une fois on a déjà le coût du patch optimal (le tableau de coût rempli).

Finalement, pour optimiser le temps et limiter les entrées/sorties vers les fichiers, nous gardons dans notre code une copie de chaque fichiers au début, faisant plus facile l'accès à l'information de chaque fichier.

## 2 Analyse du coût théorique

**Justification :** Nous allons calculer le coût de chaque fonction dans notre code, tenant en compte que il y en a quelques-unes qui s'appellent entre elles. Le patch comprend une première partie dans laquelle tous les variables globales qui ont besoin sont initialisés. Et après avec l'appel de la fonction le chemin de coût minimum est calculé. On ne décrira pas les instructions dont le coût n'est pas significative ou celles avec un coût fixe.

### 2.1 Nombre d'opérations en pire cas :

**Justification :**

- Décompte de lignes :  $n + m$   
Allocation des arrays de 1 seule dimension qui représentent des matrices
- F1 - Pour les lignes du "file source" (textes)
- F2 - Pour les lignes du "file target" (textes)
- instruction - Pour les lignes des différentes opérations (textes)
- optimal - Pour les lignes des différents coûts de opération (valeurs)
- patch - Pour les lignes du chemin optimal (textes)  
Asignation des contenues
- F1 :  $n_1$  (les lignes sont copiées du fichier source au tableau)
- F2 :  $n_2$  (les lignes sont copiées du fichier target au tableau)
- optimal :  $(n_1 + 1) * (n_2 + 1)$  (les coûts sont initialisés en infinito)
- Pire Cas de la fonction avec  $i = n_1$  et  $j = n_2$
- DelMul - Allocation d'un tableau pour stocker et analyser après le cas de destruction multiple
- Un cycle for de taille  $n_1$  pour charger le tableau DelMax
- Deux appels de la fonction findPatch de manière récursive, qui aura un coût de  $(n_1 + 1) * (n_2 + 1)$
- Génération d'un tableau pour stocker les strings des instructions du chemin optimal :  $n_1 + n_2$

Si on analyse le coût d'obtenir le chemin optimal, on observe qu'il dépend surtout des appels récursifs, c'est-à-dire que le coût théorique de l'algorithme est d'ordre  $O(n_1 * n_2)$ .

### 2.2 Place mémoire requise :

**Justification :** On va exprimer une approximation de la place en mémoire requise, en regardant que un char et un int, les deux occupent un octet :

- F1 :  $n_1 * \text{LINE\_TAILLE\_MAX}$  octets.

- F2 :  $n_2 * \text{LINE\_TAILLE\_MAX}$  octets.
- optimal :  $(n_1 + 1) * (n_2 + 1)$  octets.
- instruction :  $(n_1 + 1) * (n_2 + 1) * \text{LINE\_MAX\_SIZE}$  octets.
- patch :  $(n_1 + n_2) * \text{LINE\_MAX\_SIZE}$  octets.

Mais après les différentes appels de malloc pour obtenir les tableaux DelMul ce fait chaque fois jusqu'à l'index de l'appel, cela ajoute un coût de mémoire de  $n_1!$  octets.

### 2.3 Nombre de défauts de cache sur le modèle CO :

**Justification :** Nous voyons que les défauts de cache. Nous croyons que cela est parce que nous utilisons la méthode récursive et cela pose un problème d'allocation dans la mémoire cache. Chaque fois que le code a besoin d'un tableau, il faut le charger dans la mémoire cache, et comme la mémoire cache n'est pas très grande, il faut l'actualiser plusieurs fois.

Si l'on voit le cas de l'allocation des fichiers dans ses tableaux, on voit la cause la plus grande des défauts de cache. Tout d'abord, l'ordinateur utilisé pour les tests a une mémoire cache de 4 Mo. Nous supposons que pour la mémoire, le nombre de lignes est le même que la longueur de la ligne dedans. Dans ce cas, le nombre de lignes serait environ 1700 et la taille de chaque ligne serait d'environ 1700 octets. De plus, les lignes bien du fichier d'entrée ou bien du fichier de sortie fait au maximum 100 octets selon notre modèle. En fin de compte, on utilise une ligne de 1700 octets juste pour allouer une ligne de 100 octets, ce qui n'est pas du tout optimal, générant plus des défauts de cache.

## 3 Compte rendu d'expérimentation

### 3.1 Conditions expérimentales

#### 3.1.1 Description synthétique de la machine :

Pour les tests, on a laissé l'ordinateur sans autres exécutions pour réduire le temps d'exécution. L'ordinateur a les spécifications suivantes :

OS : Ubuntu 14.04.2 LTS  
 CPU : Intel(R) Core(TM) i7-3537U CPU @ 2.00GHz  
 GPU : NVIDIA GF117M [GeForce 610M/710M/820M / GT 620M/625M/630M/720M]  
 RAM : 5855 MiB  
 CACHE : 4096 KB

Information trouvée avec les commandes :

```
cat /proc/cpuinfo
lshw
lsb_release -a
```

#### 3.1.2 Méthode utilisée pour les mesures de temps :

Pour le temps, nous avons utilisé l'outil *temps* d'Ubuntu, en disant :

### 3.2 Mesures expérimentales

Il faut dire que nous n'avons pas pu faire l'exécution du benchmark 6, parce qu'il y avait un error pour la taille du matrice DelMax, générant une exception dans l'exécution.

### 3.3 Analyse des résultats expérimentaux

À la fin, les résultats sont un peu trop grandes, tenant en compte le nombre de défauts de cache que l'on fait. C'était impossible de trouver un moyen de comparer expérimentalement la théorie avec la pratique, parce que la matrice DelMul fait une appel a malloc, et donc le temps d'exécution n'ont pas permis réaliser la preuve.

	coût du patch	temps min	temps max	temps moyen
benchmark1	2486	5.989	7.342	6.634s
benchmark2	3093	58.572s	1m7.256s	1m1.299s
benchmark3	786	3m0.393s	3m16.467s	3m6.636
benchmark4	1638	6m.52.723s	7m39.563s	7m24.734s
benchmark5	7376	17m28.635s	18m18.974s	17m51.708s

FIGURE 1 – Mesures des temps minimum, maximum et moyen de 5 exécutions pour les 5 benchmarks.

#### 4 Question : et si le coût d'un patch était sa taille en octets ?

Pour ce cas, il faudrait changer les coûts de chaque opération, et les relier avec les octets du patch :

En cas d'une addition, la première ligne qui a un char avec le "+" et le numéro de la ligne. Et la deuxième ligne s'agit de la ligne à ajouter, ça veut dire (2 + la quantité de chiffre du nombre de caractères à ajouter). Son coût est le nombre de caractères de cette ligne.

Le cas de la substitution c'est similaire.

Pour une suppression multiple le coût est similaire mais il faut tenir en compte les chiffres des numéros de la ligne d'origine et de la quantité à supprimer.

En cas d'une suppression simple, il y a un coût de 2 + les chiffres du numéro de la ligne.

Finalement, la copie n'a pas de coût parce qu'il n'ajoute pas des lignes au patch.