



**HAL**  
open science

# Evaluation des parametres du p-hachage dynamique virtuel

G. Levy

► **To cite this version:**

G. Levy. Evaluation des parametres du p-hachage dynamique virtuel. RR-0740, INRIA. 1987. inria-00075812

**HAL Id: inria-00075812**

**<https://inria.hal.science/inria-00075812>**

Submitted on 24 May 2006

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# INRIA

UNITÉ DE RECHERCHE  
INRIA-ROCQUENCOURT

Institut National  
de Recherche  
en Informatique  
et en Automatique

Domaine de Voluceau  
Rocquencourt

BP 105

78153 Le Chesnay Cedex

France

Tel. (1) 39 63 55 11

## Rapports de Recherche

N° 740

### EVALUATION DES PARAMETRES DU p-HACHAGE DYNAMIQUE VIRTUEL

Gérard LEVY

OCTOBRE 1987

## EVALUATION DES PARAMETRES DU p-HACHAGE DYNAMIQUE VIRTUEL

Gérard Lévy\*

Les méthodes de hachage dynamique (virtuel) (HDV), ont suscité de nombreuses études qui ont permis, dans le cas  $p=2$ , d'apprécier leurs performances. La justification mathématique de ces résultats a été donnée, toujours pour  $p=2$ , par [1].

L'objet du présent article est de calculer exhaustivement, pour  $p>2$ , les paramètres du p-HDV, et d'en étudier le comportement asymptotique. Notre démarche s'inscrit dans le prolongement des travaux de [3], où l'examen du cas général est déjà amorcé.

**Mots clés** : Hachage dynamique, indexation, performances asymptotique, transformée de Mellin.

\* : Université Paris IX-Dauphine et projet SESAME-INRIA

## EVALUATION OF PARAMETERS OF DYNAMIC VIRTUAL $p$ -HASHING

Gérard LEVY \*

Dynamic (virtual) hashing methods (HDV) has led to many studies that evaluated the performance for  $p = 2$ . Mathematical justification of these results for  $p = 2$  was given in [1].

In this article we calculate in exhaustive manner the parameters of  $p$ -HDV for  $p > 2$  and we study the asymptotic behavior. Our work extends the one presented in [3], where the study of the general case was initiated.

**Key words** : Dynamic, virtual hashing, index, asymptotic performance, Mellin transform.

\* : University of Paris IX-Dauphine and SESAME Project at INRIA

## TABLE DES MATIERES

### INTRODUCTION

#### 1. RAPPELS

1.1 Hypothèse de "Bernouilli finie"

1.2 Hypothèse de "Bernouilli infinie"

1.3 Hypothèse de Poisson

#### 2. RECURRENCES VERIFIEES PAR LES PARAMETRES DU P-HACHAGE

2.1 Nombre de sommets de l'index

2.2 Coût moyen d'insertion d'une clé

2.3 Coût moyen de suppression d'une clé

2.4 Distribution des feuilles

2.5 Longueur du cheminement externe

2.6 Profondeur de l'index

#### 3. FONCTIONS GENERATRICES DES PARAMETRES

#### 4. ANALYSE ASYMPTOTIQUE DES PERFORMANCES

#### 5. ANNEXE

#### 6. CONCLUSION

#### 7. BIBLIOGRAPHIE

EVALUATION DES PERFORMANCES DU HACHAGE DYNAMIQUE VIRTUEL

INTRODUCTION

Nous allons généraliser au cas des arbres p-homogènes les résultats qui ont été obtenus pour le hachage dynamique dans le cas des arbres binaires. Outre leur intérêt mathématique, ces résultats devraient permettre de juger de l'utilité ou non de se servir d'arbres de degré supérieur à 2. Les outils et méthodes que nous employerons sont pour la plupart empruntés à [1], [2], [3]. Nous avons préféré opter pour cette approche car, en dépit de l'arsenal mathématique qu'elle nécessite, elle fournit un cadre d'étude unitaire et systématique, situation assez rare en analyse combinatoire.

Après quelques rappels mathématiques, nous donnerons les équations de récurrence des différents paramètres du hachage. Nous résoudrons ensuite ces équations. Enfin nous étudierons le comportement asymptotique des paramètres. Compte tenu de la nécessaire brièveté des rappels, le lecteur se rapportera utilement aux travaux cités en références.

I-RAPPELS

Soient p et b deux entiers positifs, p supérieur ou égal à 2. Tous les arbres dont il sera question ici sont des arbres p-homogènes aux feuilles desquels sont associées des boîtes de capacité b, à l'intérieur desquelles on peut ranger des clés. Suivant le cas ces clés sont des mots de longueur fixe s ou de longueur infinie de l'alphabet  $L = \{0, 1, \dots, p-1\} = [0..p-1]$ .

Soit  $L^s$  l'ensemble des mots de longueur s de L, et  $L^\infty$  celui des mots de longueur infinie. On désigne par  $P^{(s)}$  l'ensemble des parties w de  $L^s$  et par  $P_n^{(s)}$  celui des parties w de cardinal  $|w| = n$ , n entier naturel. (Notations similaires pour les mots infinis.)

Pour tout ensemble  $w$  de clés de  $L^s$  et tout élément  $i$  de  $L$ , on appelle  $w/i$  l'ensemble des mots  $m$  de  $L^{s-1}$  tels que  $im$  soit une clé de  $w$ .

La correspondance entre  $w$  et  $(w/0, w/1, \dots, w/p-1)$  est une bijection entre  $P^{(s)}$  et  $P^{(s-1)} \times P^{(s-1)} \times \dots \times P^{(s-1)}$ . De plus on a les relations

$$w = \bigcup_{i \in L} \{i\} \times (w/i) \quad \text{et} \quad |w| = \sum_{i \in L} |w/i| .$$

Les méthodes de hache-codage dynamique associent à tout ensemble  $w$  de clés un  $p$ -arbre de la manière suivante:

- si  $w$  contient au plus  $b$  clés il lui correspond un arbre qui n'a qu'une seule feuille dont la boîte associée renferme ces clés.
- sinon on éclate la feuille associée à  $w$  en  $p$  sous-arbres associés aux différents ensembles  $\{i\} \times (w/i)$  sur lesquels on réitère la procédure d'éclatement jusqu'à n'avoir plus que des ensembles d'au plus  $b$  clés.

Les techniques de codage précédentes font correspondre à chaque  $w$  un certain nombre de paramètres tels que la taille de l'arbre, le taux d'occupation des feuilles, le coût moyen d'une insertion ou d'une suppression de clé, etc... Ces paramètres sont intéressants par leurs valeurs moyennes, sous certaines hypothèses statistiques à préciser. Les hypothèses le plus souvent étudiées, parce que les plus plausibles, sont:

BF("Bernouilli-Finie"): les clés sont de longueur fixe  $s$  finie, la taille du "fichier"  $w$  est figée à  $n$ , et tous les  $w$  de même cardinal  $n$  sont équiprobables;

BI("Bernouilli Infinie"): mêmes hypothèses que ci-dessus, sauf que les clés sont de longueur infinie;

PS("Poisson"): mêmes hypothèses que pour BI, sauf que la taille  $n$  du fichier  $w$  fluctue en suivant une Loi de Poisson de paramètre  $\nu$ .

I-1 Hypothèse BF

On définit un paramètre  $a^{(s)}$  comme une application de  $P^{(s)}$  dans  $R$ , c-à-d un moyen d'associer à tout ensemble  $w$  de mots de  $L^s$  un réel  $a^{(s)}(w)$ .

Pour tout entier  $n$  compris entre 0 et  $p^s$ , on pose

$$a_n^{(s)} = \text{somme des } a^{(s)}(w) \text{ pour } w \text{ appartenant à } P_n^{(s)} .$$

La fonction génératrice  $a^{(s)}$  est définie par

$$a^{(s)}(x) = \sum_{w \in P^{(s)}} a^{(s)}(w) x^{|w|} = \sum_{n=0}^{p^s} \left( \sum_{w \in P_n^{(s)}} a^{(s)}(w) \right) x^n = \sum_{n=0}^{p^s} a_n^{(s)} x^n .$$

I-1-1 Propriétés des fonctions génératrices:

1) si  $a^{(s)}(w) = b^{(s)}(w) + c^{(s)}(w)$ , alors  $a^{(s)}(x) = b^{(s)}(x) + c^{(s)}(x)$  ;

2) si  $a^{(s)}(w) = \lambda b^{(s)}(w)$  alors  $a^{(s)}(x) = \lambda b^{(s)}(x)$  ;

3) si  $a^{(s)}(w) = \prod_{i \in L} b_i^{(s-1)}(w/i)$  alors  $a^{(s)}(x) = \prod_{i \in L} b_i^{(s-1)}(x)$  .

I-1-2 Valuations particulières :

1) si  $a^{(s)}(w) = 1$  alors  $a_n^{(s)} = |P_n^{(s)}| = C_p^n = \binom{p}{n}$  et  $a^{(s)}(x) = (1+x)^{p^s}$  ;

2) si  $a^{(s)}(w) = \delta_{|w|,q}$  ( $=1$  si  $|w|=q$ ,  $=0$  sinon), alors  $a_q^{(s)} = \binom{p}{q}$ , et

$$a^{(s)}(x) = \binom{p}{q} x^q ;$$

3) si  $a^{(s)}(w) = |w|$ , alors  $a_n^{(s)} = n \binom{p}{n}$  et  $a^{(s)}(x) = \sum_{n=0}^{p^s} n \binom{p}{n} x^n = p^s x (1+x)^{p^s-1}$

4) si  $a^{(s)}(w) = b_i^{(s-1)}(w/i)$  alors  $a^{(s)}(w) = b_i^{(s-1)}(w/i) \prod_{j \in L - \{i\}} a^{(s-1)}(w/j)$ , d'où

$$a^{(s)}(x) = b_i^{(s-1)}(x) (1+x)^{(p-1)p^{s-1}} = b_i^{(s-1)}(x) (1+x)^{p^s - p^{s-1}}.$$

5) récurrences: si  $a^{(s)}(w) = \sum_{i \in L} a^{(s-1)}(w/i) + b^{(s)}(w)$ , alors

$$a^{(s)}(x) = p(1+x)^{(p-1)p^{s-1}} \cdot a^{(s-1)}(x) + b^{(s)}(x).$$

Or si  $a^{(s)}(x) = c_s(x) \cdot a^{(s-1)}(x) + d_s(x)$ , avec  $a^{(0)}(x) = c_0(x)$ , alors on a

$$a^{(s)}(x) = \sum_{j=0}^s d_j(x) \prod_{k=j+1}^s c_k(x).$$

$$\text{Donc } a^{(s)}(x) = \sum_{j=0}^s d_j(x) \cdot p^{s-j} \cdot (1+x)^{p^s - p^j}.$$

### I-2 Hypothèse BI

Un paramètre  $a^{(\infty)}$  est une application de  $P^{(\infty)}$  dans  $\mathbb{R}$ .

Pour tout entier naturel  $n$ , on définit

$a_n^{(\infty)} = E\{a^{(\infty)}(w) \mid w \in P_n^{(\infty)}\}$ , et la série génératrice associée

$$a^{(\infty)}(x) = \sum_{n=0}^{\infty} a_n^{(\infty)} (x^n/n!)$$

Sous les hypothèses d'équiprobabilité que nous avons prises ici, on a

$$a_n^{(\infty)} = \lim_{s \rightarrow \infty} a_n^{(s)} / \binom{p^s}{n} = \lim_{s \rightarrow \infty} a_n^{(s)} / ((p^s)^n/n!);$$

$$\text{Donc } a_n^{(\infty)}/n! = \lim_{s \rightarrow \infty} a_n^{(s)} / (p^s)^n, \text{ et } a^{(\infty)}(x) = \lim_{s \rightarrow \infty} a^{(s)}(x/p^s).$$

Nous allons appliquer à l'étude de l'hypothèse BI le même ordre et les mêmes notations qu'à l'étude de BF, en prenant soin toutefois de remplacer les indices supérieurs  $s$  et  $s-1$  par l'indice  $\infty$ .

I-2-1 Propriétés des fonctions génératrices:

1), 2), propriétés de linéarité sont vérifiées;

$$3) \text{ si } a^{(\infty)}\{w\} = \prod_{i \in L} b_i^{(\infty)}\{w/i\} \text{ alors } a^{(\infty)}\{x\} = \prod_{i \in L} b_i^{(\infty)}\{x/p\}.$$

I-2-2 Valuations particulières:

$$1) \text{ si } a^{(\infty)}\{w\} = 1 \text{ alors } a_n^{(\infty)} = 1 \text{ et } a^{(\infty)}\{x\} = \sum_{n=0}^{\infty} 1(x^n/n!) = e^x ;$$

$$2) \text{ si } a^{(\infty)}\{w\} = \delta_{|w|,q} \text{ alors } a^{(\infty)}\{x\} = x^q/q! ;$$

$$3) \text{ si } a^{(\infty)}\{w\} = |w| \text{ alors } a^{(\infty)}\{x\} = \sum_{n=0}^{\infty} n(x^n/n!) = xe^x ;$$

$$4) \text{ si } a^{(\infty)}\{w\} = b_i^{(\infty)}\{w/i\} \text{ alors } a^{(\infty)}\{x\} = b_i^{(\infty)}\{x/p\} \cdot e^{\frac{p-1}{p}x}$$

$$5) \text{ récurrences: si } a^{(\infty)}\{w\} = \sum_{i \in L} a^{(\infty)}\{w/i\} + b^{(\infty)}\{w\}, \text{ alors}$$

$$a^{(\infty)}\{x\} = p \cdot e^{\frac{p-1}{p}x} \cdot a^{(\infty)}\{x/p\} + b^{(\infty)}\{x\}.$$

Or si  $a^{(\infty)}\{x\} = c(x) \cdot a^{(\infty)}\{x/p\} + d(x)$ , on a pour tout  $k$  naturel

$$a^{(\infty)}\{x\} = \sum_{j=0}^k d(x/p^j) \cdot \prod_{i=0}^{j-1} c(x/p^i) + a^{(\infty)}\{x/p^{k+1}\} \cdot \prod_{i=0}^k c(x/p^i).$$

Avec  $c(x) = p \cdot \exp(\frac{p-1}{p}x)$ , on obtient, si les produits convergent,

$$a^{(\infty)}(x) = e^x \cdot \sum_{j=0}^{\infty} p^j \cdot b^{(\infty)}\{x/p^j\} \cdot e^{-(x/p^j)} .$$

I-3 Hypothèse P

Les clés sont de longueur infinie et la taille du fichier est une variable aléatoire  $N$  qui suit une Loi de Poisson de paramètre  $\nu$ . Donc

$\text{Prob}(N=n) = e^{-\nu} \cdot (\nu^n/n!)$ , et la valeur moyenne  $a^{(P)}$  du paramètre  $a_n^{(\infty)}$  est

$$a^{(P)} = \sum_{n=0}^{\infty} a_n^{(\infty)} \cdot e^{-\nu} \cdot (\nu^n/n!) = e^{-\nu} \cdot \sum_{n=0}^{\infty} a_n^{(\infty)} (\nu^n/n!) = e^{-\nu} \cdot a^{(\infty)}(\nu) .$$

II-RECURRENCES VERIFIEES PAR LES PARAMETRES DU HACHE-CODAGE.

Pour s fini ou non, nous définissons les fonctions caractéristiques  $\chi$  à valeurs 0 ou 1 suivant qu'une proposition concernant  $w$  de  $P^{(s)}$  est fausse ou vraie. Ainsi

$$\chi(|w| \leq b) = 1 \text{ si } |w| \leq b, \text{ et } 0 \text{ sinon.}$$

On a évidemment

$$\delta_{|w|,j} = \chi(|w|=j) \quad , \quad \chi(|w| \leq b) = \sum_{j=0}^b \delta_{|w|,j} = 1 - \chi(|w| > b) .$$

L'arbre associé à l'ensemble  $w$  de clés s'appelle l'index, et ce sont ses paramètres que l'on étudie.

II-1 Nombre de sommets internes de l'index.

Ce nombre est égal à 0 si  $w$  n' a pas plus de  $b$  clés. Sinon il est égal à la somme des nombres de sommets internes des index des différents  $w/i$  plus 1. Donc

$$N^{(s)}(w) = (\sum_{i \in L} N^{(s-1)}(w/i) + 1) \chi(|w| > b) , \text{ soit}$$

$$N^{(s)}(w) = \sum_{i \in L} N^{(s-1)}(w/i) + 1 - \chi(|w| \leq b) .$$

Le nombre de feuilles de l'index est lié au nombre de sommets internes par l'égalité déjà rencontrée

$$N^{(s)}(w) = (F^{(s)}(w) - 1) / (p - 1), \text{ d'où } F^{(s)}(w) = 1 + (p - 1) \cdot N^{(s)}(w).$$

II-2 Coût moyen d'insertion d'une clé dans un fichier w de n clés.

Soit  $I^{(s)}(w)$  ce coût. Il est égal à la moyenne des coûts d'insertion de la nouvelle clé dans les différents w/i, si w contient plus de b clés. Sinon il vaut 1 si w en contient b exactement, puisqu'alors un éclatement est nécessaire. D'où la récurrence

$$I^{(s)}(w) = (1/p) \cdot \left( \sum_{i \in L} I^{(s-1)}(w/i) \right) + \delta_{|w|, b}.$$

Une autre manière d'évaluer ce nombre consiste à dire que pour chaque n  $I_n^{(s)}$  est égal au nombre d'éclatements résultant des insertions, donc à la différence  $N_{n+1}^{(s)} - N_n^{(s)}$  des nombres de sommets internes des index des fichiers à n+1 et à n clés, puisque chaque éclatement a pour effet d'ajouter un sommet interne nouveau à l'index. D'où

$$I_n^{(s)} = N_{n+1}^{(s)} - N_n^{(s)}.$$

II-3 coût moyen de suppression .

Avec les mêmes notations que çï-dessus, w ayant n clés, la suppression étant l'opération inverse de l'insertion, on peut considérer qu'en moyenne le nombre  $S^{(s)}(w)$  ou  $S_n^{(s)}$  de suppressions dans un fichier de n clés est égal au nombre d'insertions dans un fichier de n-1 clés. D'où

$$S_n^{(s)} = I_{n-1}^{(s)}.$$

II-4-Distribution des feuilles.

Soit  $F_q$  le nombre des feuilles de l'index qui contiennent exactement  $q$  feuilles,  $q$  entier compris entre 0 et  $b$ .

Si  $w$  a plus de  $b$  clés ce nombre est égal à la somme des nombres  $F_q$  correspondants aux différents  $w/i$ . Sinon il vaut 1 si  $w$  contient  $q$  clés, Et 0 sinon. D'où

$$F_q^{(s)}(w) = \left( \sum_{i \in L} F_q^{(s-1)}(w/i) \right) \cdot \chi(|w| > b) + \delta_{|w|, q}$$

$$= \sum_{i \in L} F_q^{(s-1)}(w/i) + \delta_{|w|, q} - \sum_{i \in L} F_q^{(s-1)}(w/i) \cdot \chi(|w| \leq b)$$

Or si  $w$  a au plus  $b$  clés, il en est de même de  $w/i$  et  $F_q^{(s-1)}(w/i) = 0$  ou 1. Plus précisément, ce nombre vaut 1 ssi  $w/i$  a  $q$  clés et l'ensemble des autres  $w/j, j \neq i$ , en contient au plus  $b-q$ . Donc

$$F_q^{(s-1)}(w/i) \cdot \chi(|w| \leq b) = \delta_{|w/i|, q} \cdot \chi\left(\sum_{j \in L - \{i\}} |w/j| \leq b - q\right)$$

Posons  $E_k$  = ensemble des  $w$  tels que la somme des  $|w/j|, j \neq i$ , est égale à  $k$ .

Désignons par  $R_k$  l'ensemble des suites  $(r_j)$  de  $p-1$  entiers naturels  $r_j$  telles que la somme de ces  $r_j$  soit égale à  $k$ . On a les égalités

$$\{w \mid \sum_{j \neq i} |w/j| \leq b - q\} = \bigcup_{k=0}^{b-q} E_k = \bigcup_{k=0}^{b-q} \bigcup_{(r_j) \in R_k} \{w \mid (w/j) = (r_j)\}$$

$$\chi\left(\sum_{j \neq i} |w/j| \leq b - q\right) = \sum_{k=0}^{b-q} \sum_{(r_j) \in R_k} \prod_{j \neq i} \delta_{|w/j|, r_j}$$

D'où l'expression de la récurrence relative aux  $F_q$  :

$$F_q^{(s)}(w) = \sum_{i \in L} F_q^{(s-1)}(w/i) + \delta_{|w|,q} - \sum_{i \in L} \delta_{|w/i|,q} \cdot \sum_{k=0}^{b-q} \sum_{(r_j) \in R_k} \prod_{j \neq i} \delta_{|w/j|,r_j}$$

II-5 Longueur du cheminement externe.

C'est la somme des longueurs des chemins conduisant de la racine aux feuilles de l'index. Les nombres  $L^{(s)}(w)$  satisfont la relation

$$\begin{aligned} L^{(s)}(w) &= \left( \sum_{i \in L} (L^{(s-1)}(w/i) + |w/i|) \right) \cdot \chi(|w| > b) \\ &= \sum_{i \in L} L^{(s-1)}(w/i) + |w| - \sum_{j=0}^b j \delta_{|w|,j} \end{aligned}$$

II-6 Profondeur de l'index.

sous l'hypothèse BI, soit  $P_n^{(\infty)}(k)$  la probabilité qu'un index de taille  $n$  soit de profondeur au plus égale à  $k$ . Posons

$P_k^{(\infty)}(w) = 1$  si la profondeur de  $w$  est au plus  $k$ , et 0 sinon.

Puisque  $P_k^{(\infty)}(n) = E\{ P_k^{(\infty)}(w) \mid |w| = n \}$ , on a  $P_k^{(\infty)}(n) = P_n^{(\infty)}(k)$ .

De plus  $P_{k+1}^{(\infty)}(w) = \prod_{i \in L} P_k^{(\infty)}(w/i)$ .

De ces relations, on déduit que la probabilité qu'un index de taille  $n$  soit égale à  $k$ , vaut  $P_n^{(\infty)}(k) - P_n^{(\infty)}(k-1)$ , et que la profondeur moyenne de l'index est, pour un fichier de taille  $n$ ,

$$PF(n) = \sum_{k=1}^{\infty} k (P_n^{(\infty)}(k) - P_n^{(\infty)}(k-1))$$

III-FONCTIONS GENERATRICES DES PARAMETRES.

Nous allons nous servir des outils fournis en I pour les appliquer aux récurrences de II, dans l'ordre où elles ont été obtenues, en vue d'obtenir les fonctions génératrices des différents paramètres des indices. De ces fonctions nous déduirons l'expression des paramètres, puis ultérieurement leur comportement asymptotique.

III-1 nombre de sommets internes.

Hypothèse BF: on a  $b^{(s)}(w) = 1 - \chi(|w| \leq b) = 1 - \sum_{j=0}^b \delta_{|w|, j}$ .

Donc  $b^{(s)}(x) = d_s(x) = (1+x)^p - \sum_{j=0}^b \binom{p}{j} x^j$ . ET

$$N^{(s)}(x) = \sum_{k=1}^s p^{s-k} (1+x)^{p^{s-k}} \cdot \left( \sum_{j=b+1}^p \binom{p}{j} x^j \right)$$

$$= \sum_{k=1}^s p^{s-k} \left[ \sum_{n=b+1}^p \left( \sum_{j=b+1}^n \binom{p}{j} \binom{p^{s-k}}{n-j} \right) x^n \right]$$

Il en résulte que le coefficient de  $x^n$  vaut

$$N_n^{(s)} = \sum_{k=1}^s p^{s-k} \left( \sum_{j=b+1}^n \binom{p}{j} \binom{p^{s-k}}{n-j} \right), \text{ si } n \geq b+1, \text{ et } 0 \text{ si } n \leq b.$$

Hypothèse BI:

$$b^{(\infty)}(x) = e^x - \sum_{j=0}^b (x^j/j!) = e^x - e_b(x) = \sum_{j=b+1}^{\infty} (x^j/j!).$$

$$N^{(\infty)}(x) = \sum_{k=0}^{\infty} p^k \cdot \exp(1-1/p^k) \cdot \left[ \exp(x/p^k) - \sum_{j=0}^b (x/p^k)^j/j! \right].$$

Le coefficient de  $x^n/n!$  est donné par

$$N_n^{(\infty)} = \sum_{k=0}^{\infty} p^k \left[ 1 - \sum_{j=0}^b \binom{n}{j} (1/p^k)^j (1-1/p^k)^{n-j} \right]$$

Le taux moyen de remplissage des boîtes attachées aux feuilles de l'index est défini par

$$\tau_n = n / (bF_n) = n / (b(1 + (p-1)N_n)), \text{ et } \tau = \lim_{n \rightarrow \infty} \tau_n$$

### III-Coût moyen d'insertion.

On a les égalités suivantes

$$I^{(s)}(x) = p \cdot (1/p) \cdot I^{(s-1)}(x) + \binom{p^s}{b} \cdot x^b, \text{ soit } c_s(x) = (1+x)^{p^s - p^{s-1}}, \text{ et}$$

$$d_s(x) = \binom{p^s}{b} \cdot x^b. \text{ D'où } I^{(s)}(x) = (1+x)^{p^s - p^{s-1}} I^{(s-1)}(x) + \binom{p^s}{b} \cdot x^b;$$

$$I^{(s)}(x) = x^b \cdot \sum_{j=0}^s \binom{p^j}{b} \cdot (1+x)^{p^s - p^j}$$

Ce qui permet de déterminer les  $I_n^{(s)}$ , en développant l'expression

$$I^{(s)}(x) = x^b \cdot \sum_{j=0}^s \binom{p^j}{b} \cdot \sum_{k=0}^{p^s - p^j} \binom{p^s - p^j}{k} x^k = \sum_{k=0}^{p^s - 1} \left( \sum_{j=0}^{\text{Log}_p(p^s - k)} \binom{p^j}{b} \binom{p^s - p^j}{k} \right) x^{b+k}$$

$$\text{D'où } I_n^{(s)} = \sum_{j=0}^{\text{Log}_p(p^s - n + b)} \binom{p^j}{b} \binom{p^s - p^j}{n - b}$$

Sous l'hypothèse BI on a

$$I^{(\infty)}(x) = I^{(\infty)}(x/p) \cdot \exp((1-1/p)x) + x^b/b!$$

Ce qui entraîne que

$$I^{(\infty)}(x) = \sum_{k=0}^{\infty} (1/b!) \cdot e^x \cdot e^{-(x/p^k)} = (x^b/b!) \cdot \sum_{k=0}^{\infty} (1/p^k)^b \cdot e^{x(1-1/p^k)}$$

Le coefficient de  $x^n/n!$  permet d'obtenir

$$I_n^{(\infty)} = \binom{n}{b} \cdot \sum_{k=0}^{\infty} (1/p^k)^b (1-1/p^k)^{n-b}$$

Sous l'hypothèse P de Poisson, la valeur moyenne du coût d'insertion est

$$I^{(P)}(v) = e^{-v} \cdot I^{(\infty)}(v) = (v^b/b!) \cdot \sum_{k=0}^{\infty} (1/p^k)^b \cdot e^{-v/p^k}$$

III-3 Coût moyen de suppression.

Son étude découle immédiatement de celle du coût d'insertion, puisque

$$S_n = I_{n-1}$$

III-4- Distribution des feuilles.

À la différence des autres paramètres, les seconds membres des relations de récurrence relatives à la repartition des feuilles ont des expressions assez compliquées. Il nous faudra donc en un premier lieu, calculer leurs fonctions génératrices. Nous allons poser,

$$\text{pour simplifier les notations } g_i((w/j):j \neq i) = \chi\left(\sum_{j \in L-\{i\}} |w/j| \leq b-q\right)$$

Une expression détaillée de  $g_i$  a été déjà calculée.

Hypothèse BF:

On sait que dans  $L^{(s)}$ , il correspond à  $\delta_{|w|,q}$  la fonction génératrice

$\binom{p}{q} x^q$ . Il correspond donc à  $g_i$ , compte tenu de son expression détaillée,

$$\sum_{k=0}^{b-q} \sum_{(r_j) \in R_k} \prod_{j \neq i} \binom{p^{s-1}}{r_j} x^{r_j} = \sum_{k=0}^{b-q} \sum_{(r_j) \in R_k} x^{(\sum r_j)} \prod_{j \neq i} \binom{p^{s-1}}{r_j}$$

Or pour toute suite  $(r_j)$  de  $R_k$ , la somme des  $r_j$  est égale à  $k$ . Donc

$$g_i(x) = \sum_{k=0}^{b-q} x^k \sum_{(r_j) \in R_k} \prod_{j \neq i} \binom{p^{s-1}}{r_j} = \sum_{k=0}^{b-q} \binom{p^s - p^{s-1}}{k} x^k$$

Cette égalité s'obtient aisément si on remarque que

$$p^s - p^{s-1} = (p-1)p^{s-1}, \text{ et donc que } (1+x)^{p^s - p^{s-1}} = \prod_{j \in L - \{i\}} (1+x)^{p^{s-1}} =$$

$$= \prod_{j \neq i} \left( \sum_{r_j=0}^{p^{s-1}} \binom{p^{s-1}}{r_j} x^{r_j} \right) = \sum_{k=0}^{p^s} \left( \sum_{(r_j) \in R_k} x^{(\sum r_j)} \prod_{j \neq i} \binom{p^{s-1}}{r_j} \right)$$

Donc 
$$\sum_{(r_j) \in R_k} \prod_{j \neq i} \binom{p^{s-1}}{r_j} = \binom{p^s - p^{s-1}}{k} .$$

Si on convient de poser pour toute fonction f qui admet le développement en série

$$f(x) = \sum_{n=0}^{\infty} a_n x^n, \quad \llbracket f(x) \rrbracket_i^j = \sum_{n=i}^j a_n x^n, \text{ on constate que}$$

$$g_i(x) = \llbracket (1+x)^{p^s - p^{s-1}} \rrbracket_0^{b-q} .$$

Il correspond, à présent, à  $\delta_{|w/i|, q} \cdot g_i((w/j))$  la fonction

$$\binom{p^{s-1}}{q} x^q \cdot g_i(x), \text{ et au second membre tout entier de la récurrence,}$$

$$b_q^{(s)}(x) = d_s(x) = \binom{p^s}{q} x^q - p \binom{p^{s-1}}{q} x^q \left( \sum_{k=0}^{b-q} \binom{p^s - p^{s-1}}{k} x^k \right) .$$

Les seconds membres relatifs aux feuilles vides et pleines sont

$$b_0^{(s)}(x) = 1 - p \sum_{k=0}^b \binom{p^s - p^{s-1}}{k} x^k, \text{ et } b_b^{(s)}(x) = \left( \binom{p^s}{b} - p \binom{p^{s-1}}{b} \right) x^b .$$

La fonction génératrice  $F_q^{(s)}(x)$  du nombre de feuilles qui contiennent  $q$  clés est donc

$$F_q^{(s)}(x) = \sum_{j=0}^s b_q^{(j)}(x) \cdot p^{s-j} \cdot (1+x)^{p^s - p^j}$$

$$= x^q \sum_{j=0}^s p^{s-j} \left( \binom{p^j}{q} - p \binom{p^{j-1}}{q} \right) (1+x)^{p^s - p^j} \left( \sum_{k=0}^{b-q} \binom{p^s - p^{s-1}}{k} x^k \right)$$

et le nombre  $F_{q,n}^{(s)}$  est le coefficient de  $x^n$  de  $F_q^{(s)}(x)$ .

Hypothèse BI.

Ici encore c'est la détermination de  $b_q^{(\infty)}(x)$  qui pose problème.

Du fait qu'à  $\delta_{|w/j|, r_j}$  il correspond  $(x/p)^{r_j} / r_j!$ , il correspond à

$g_i((w/j))$  la fonction

$$g_i(x) = \sum_{k=0}^{b-q} \sum_{(r_j) \in R_k} \prod_{j \neq i} (x/p)^{r_j} / r_j! = \sum_{k=0}^{b-q} \frac{(x^k/k!)}{p^k} \sum_{(r_j) \in R_k} \left( \frac{k!}{\prod_{j \neq i} r_j!} \right)$$

Or la somme des termes portant sur les  $(r_j)$  de  $R_k$  est égale à  $(p-1)^k$ .

$$D'où \quad g_i(x) = \sum_{k=0}^{b-q} \left( \frac{p-1}{p} x \right)^k / (k!) = e_{b-q} \left( \frac{p-1}{p} x \right) = \left[ \exp \left( \frac{p-1}{p} x \right) \right]_0^{b-q}.$$

A titre indicatif, on remarquera qu'une manière de calculer le coefficient de  $x^k / (k!)$  dans  $g_i(x)$  est d'évaluer de deux façons différentes ce coefficient dans

$$\left( \frac{x}{ep} \right)^{p-1} = \prod_{j \neq i} e^{\frac{x}{p}} = \prod_{j \neq i} \sum_{r_j=0}^{\infty} (x/p)^{r_j} / (r_j!) = \sum_{k=0}^{\infty} \frac{(p-1)^k}{p^k} x^k / (k!).$$

Au second membre tout entier de la récurrence correspond enfin

$$b_q^{(\infty)}(x) = (x^q/q!) - p((x/p)^q/q!) \cdot e_{b-q}^{p-1}(\frac{x}{p}) = (x^q/q!)[1 - p^{1-q} \cdot e_{b-q}(\frac{x}{p})]$$

$$D'où F_q^{(\infty)}(x) = e^x \cdot \sum_{j=0}^{\infty} p^j \cdot b_q^{(\infty)}(x/p^j) \cdot e^{-(x/p^j)}$$

$$= e^x \cdot (x^q/q!) \cdot \sum_{j=0}^{\infty} p^{(1-q)j} [1 - p^{1-q} \cdot e_{b-q}(\frac{x}{p^j})] \cdot e^{-(x/p^j)}$$

Le calcul des coefficients de  $x^n$  dans les séries génératrices conduit aux résultats suivants:

sous des hypothèses de Bernouilli, pour des fichiers  $w$  de  $n$  clés, le nombre de feuilles de l'index qui contiennent  $q$  clés est

$$F_{q,n}^{(s)} = \delta_{n,q} + \sum_{k=0}^s p^{(s-k)} \left[ \binom{p}{q} \binom{p-s-k}{n-q} - p \binom{p-k-1}{q} \sum_{i=0}^{b-q} \binom{p-k-1}{i} \binom{p-s-k}{n-q-i} \right];$$

$$F_{q,n}^{(\infty)} = \delta_{n,q} + \binom{n}{q} \sum_{k=0}^{\infty} p^{(1-q)k} \left[ (1-p^{-k})^{n-q} - p^{1-q} \sum_{i=0}^{b-q} \binom{n-q}{i} \left(\frac{p-1}{p}\right)^{-k} (1-p^{-k})^{n-q-i} \right]$$

### III-5-Longueur du cheminement externe.

Au second membre  $|w| - \sum_{j=0}^b j \delta_{|w|,j}$  correspond dans le cas BF

$$b^{(s)}(x) = p^s x(1+x)^{p^s-1} - \sum_{j=0}^b j \binom{p^s}{j} x^j = \sum_{j=b+1}^{p^s} j \binom{p^s}{j} x^j; \text{ d'où}$$

$$L^{(s)}(x) = \sum_{k=0}^s p^{s-k} b^{(k)}(x) \cdot (1+x)^{p^s-p^k}$$

$$= \sum_{k=0}^s p^{s-k} \cdot \left( \sum_{i=b+1}^{p^s} i \binom{p^s}{i} x^i \right) \cdot \left( \sum_{j=0}^{p^s-p^k} \binom{p^s-p^k}{j} x^j \right)$$

Dans le cas BI, il correspond au second membre la fonction

$$b^{(\infty)}(x) = xe^x - \sum_{j=0}^{b-1} j(x^j/j!) = x(e^x - e_{b-1}(x)). \text{ D'où}$$

$$L^{(\infty)}(x) = e^x \sum_{k=0}^{\infty} p^k \cdot b^{(\infty)}(x/p^k) \cdot e^{-(x/p^k)}$$

$$= xe^x(1 - e^{-x} e_{b-1}(x)) + xe^x \sum_{k=1}^{\infty} [1 - e^{-(x/p^k)} (1 + \sum_{j=1}^{b-1} \frac{1}{j!} \frac{x^j}{(p^k)^j})]$$

L'extraction des coefficients de  $x^n$  dans ces fonctions génératrices donne les valeurs du paramètre  $L_n$  sous les deux hypothèses BF et BI.

III-6- Profondeur de l'index.

Sous l'hypothèse BI, on a vu que  $P_{k+1}(w) = \prod_{i \in L} P_k(w/i)$ .

De plus on a  $P_0(w)=1$  si et seulement si  $w$  n'a pas plus de  $b$  clés. Donc

$$P_k(x) = (P_{k-1}(x/p))^p = \dots = (P_0(x/p^k))^p, \text{ avec}$$

$$P_0(x) = \sum_{j=0}^b x^j/j! = e_b(x) \text{ . Par conséquent}$$

$$P_k(x) = (e_b(x/p^k))^p \text{ .}$$

Il s'ensuit que la probabilité que l'index d'un fichier de  $n$  clés soit de profondeur inférieure ou égale à  $k$  est :

$$P_n(k) = P_k(n) = n! [x^n] (e_b(x/p^k))^p .$$

Dans le cas poissonien, on obtient la probabilité que la profondeur du fichier soit inférieure ou égale à  $k$

$$P_v(k) = e^{-v} . P_k(v) = (f_b(vp^{-k}))^p , \text{ avec } f_b(x) = e^{-x} . e_b(x) .$$

Quant à la profondeur moyenne de l'index, elle est donnée par

$$PF(n) = \sum_{k=1}^{\infty} k(P_n(k) - P_n(k-1)) = \sum_{k=0}^{\infty} (1 - P_n(k)) ,$$

ainsi qu'on peut s'en rendre compte par transformation d'Abel.

Sous l'hypothèse de Poisson, on a un résultat similaire en passant aux valeurs moyennes.

IV- ANALYSE ASYMPTOTIQUE DES PERFORMANCES

Nous allons nous intéresser au comportement asymptotique de certains des paramètres du hachage dynamique virtuel sous la seule hypothèse Poissonienne, lorsque la taille moyenne  $v$  du fichier  $w$  devient infinie. On constate, comme on le verra un peu plus loin, que l'expression de plusieurs de ces paramètres utilise des séries infinies d'un type particulier qui ont déjà été étudiées pour  $p=2$ , et dont l'étude se généralise immédiatement pour  $p$  plus grand que 2.

Définition:

Posons  $S_j(x) = \sum_k p^{(1-j)k} \cdot e^{-xp^{-k}}$ ,  $p, j$  entiers supérieurs à 1,  $x$  positif.

Lemme Asymptotiquement en  $x$ , on a

$$S_j(x) = x^{1-j} \cdot \frac{\Gamma(j-1)}{\text{Ln}(p)} \left( 1 + \sum_{l \neq 0} e^{2i\pi \text{Log}_p x} \frac{\Gamma(j-1+2il\pi/\text{Ln}(p))}{\Gamma(j-1)} + O(x^{-m}) \right),$$

où  $m$  est un réel positif quelconque.

La preuve est identique à celle de [ ] où l'on remplace 2 par  $p$ .

La somme qui apparaît dans les  $S_j(x)$  est une fonction périodique de  $\text{Log}_p x$  de période 1, de moyenne nulle, notée  $N_j(\text{Log}_p x)$ .

IV- Nombre de sommets internes de l'index .

On a vu que le second membre, pour  $s$  infini, est  $b(x) = \sum_{j=b+1}^{\infty} x^j / j!$ .

$$\text{Donc } N^{(\infty)}(x) = e^x \sum_{k=0}^{\infty} p^k \cdot b(x/p^k) \cdot e^{-x/p^k} = e^x \sum_{k=0}^{\infty} p^k \left( \sum_{j=b+1}^{\infty} \frac{x^j}{(p^k)^j} \right) e^{-x/p^k}$$

$$= e^x \sum_{j=b+1}^{\infty} \frac{x^j}{j!} \left( \sum_{k=0}^{\infty} p^{(1-j)k} \cdot e^{-x/p^k} \right) = e^x \sum_{j=b+1}^{\infty} \frac{x^j}{j!} S_j(x)$$

Or  $N^{(P)}(v) = e^{-v} \cdot N^{(\infty)}(v)$  . Donc

$$N^{(P)}\{v\} = \sum_{j=b+1}^{\infty} (v^j/j!) \cdot s_j(v) .$$

Compte tenu du lemme , on a , lorsque le paramètre devient infini :

$$N^{(P)}\{v\} = \sum_{j=b+1}^{\infty} (1/j!) \cdot v^j \cdot v^{1-j} \cdot \frac{\Gamma(j-1)}{\text{Ln}(p)} \cdot (1 + N_j(\text{Log}_p v) + O(v^{-m})) .$$

$$= (v/\text{Ln}(p)) \cdot \sum_{j=b+1}^{\infty} \frac{1}{j(j-1)} (1 + N_j(\text{Log}_p v) + O(v^{-m})) .$$

Or la série de terme  $1/j(j-1)$ ,  $j$  allant de  $b+1$  à l'infini, est convergente et a pour somme  $1/b$  . Par conséquent on a le

#### Théorème

Le nombre moyen de sommets internes de l'index d'un fichier dont la taille est une variable aléatoire de Poisson de paramètre  $v$ , est asymptotiquement en  $v$  :

$$N^{(P)}\{v\} = \frac{v}{b \text{Ln}(p)} (1 + N_S(\text{Log}_p v) + O(v^{-m})) ,$$

où  $m$  est un réel positif quelconque , et  $N_S$  une fonction périodique de période  $1$ , de moyenne nulle de  $\text{Log}_p v$  , et dont les coefficients de

Fourier sont

$$c_1 = \sum_{j=b+1}^{\infty} \frac{b}{j(j-1)} \cdot \frac{\Gamma(j-1+2i1\pi/\text{Ln}(p))}{\Gamma(j-1)} , \text{ pour } 1 \neq 0 .$$

Puisque le nombre total de feuilles de l'index  $F^{(P)}$  est lié au nombre de sommets internes par  $F = 1 + (p-1)N$ , on en déduit facilement une expression de  $F$ . De cette dernière découle le taux de remplissage asymptotique

Corollaire:

Sous les mêmes hypothèses que précédemment, le taux de remplissage asymptotique des boîtes attachées aux feuilles de l'index est

$$\tau = \frac{\text{Ln}(p)}{p-1}$$

Remarque /pour  $p=2$ , on retrouve la valeur connue  $\tau = \text{Ln}(2)$ . De plus on constate que le taux de remplissage moyen est une fonction décroissante de  $p$ .

Statistiques sur les feuilles.

Nous allons effectuer la même démarche pour calculer le comportement asymptotique du nombre  $F_q^{(P)}(v)$ . On a

$$\begin{aligned} F_q^{(P)}(x) &= \sum_k p^k \cdot b_q(x/p^k) \cdot e^{-x/p^k} \\ &= \sum_k p^k (1/q!) \cdot (x/p^k)^q \cdot [1 - p^{1-q} \sum_{j=0}^{b-q} (1/j!) (\frac{p-1}{p})^j \cdot (x/p^k)^j] \cdot e^{-x/p^k} \\ &= (x^q/q!) [S_q(x) - p^{1-q} \sum_{j=0}^{b-q} (1/j!) (\frac{p-1}{p})^j x^j S_{j+q}(x) ]. \end{aligned}$$

Pour  $q$  au moins égal à 2, on peut utiliser le lemme. on obtient

Théorème:

asymptotiquement en  $v$ , le nombre moyen de feuilles de l'index qui contiennent  $q$  clés est

$$F_q^{(P)}(v) = \frac{v}{q! \text{Ln}(p)} \left[ \Gamma(q-1) - p^{1-q} \sum_{j=0}^{b-q} \frac{1}{j!} \left(\frac{p-1}{p}\right)^j \Gamma(j+q-1) \right] \left[ 1 + F_S(\text{Log}_p v) + O(v^{-m}) \right]$$

où  $F_S(\text{Log}_p v)$  est une fonction périodique de période 1, de moyenne nulle, dont les coefficients de Fourier se calculent comme ceux de  $N^{(P)}(v)$ , et où  $m$  est un réel positif quelconque. (Résultat vrai pour  $q \geq 2$ )

Corollaire

Asymptotiquement, le nombre moyen de feuilles qui contiennent  $q$  clés est

$$F_q^{(P)}(v) \sim \frac{v}{q(q-1)\text{Ln}(p)} \left[ 1 - p^{1-q} \sum_{j=0}^{b-q} \binom{j+q-2}{j} \left(\frac{p-1}{p}\right)^j \right], \text{ pour } q \geq 2.$$

Preuve : elle se fonde sur le fait que  $\Gamma(n) = (n-1)!$ .

En particulier si  $q=b$ , le nombre moyen de feuilles pleines est donné par

$$F_b^{(P)}(v) \sim \frac{v}{b(b-1)\text{Ln}(p)} (1 - p^{1-b}).$$

Remarque: ici encore on constate que le nombre de feuilles pleines est fonction décroissante de  $p$ . Ce qui était, bien sûr, prévisible.

Cas des feuilles vides ( $q=0$ )

Ainsi qu'on l'a vu, les résultats précédents ne sont valables que si  $q$  est supérieur à 1. Cela vient du fait que les séries  $S_j(x)$  divergent pour  $q=0$  ou 1. Or l'étude du nombre de feuilles vides est un indicateur intéressant dans le problème de hachage, puisqu'il est égal au nombre de cases inoccupées. Pour déterminer ce nombre nous allons suivre la

démarche de [1] qui consiste à déterminer le nombre de sommets internes de l'index qui ont au moins un fils qui contient q clés exactement.

Ces nombres  $M_q^{(\infty)}(w)$  sont liés, si le cardinal de w est supérieur à b, par

$$M_q(w) = \sum_{i \in L} M_q(w/i) + \sum_{i \in L} \delta_{|w/i|, q} \cdot \chi(\sum_{j \neq i} |w/j| > b-q) .$$

De plus  $M_q(w) = \delta_{|w|, q}$ , si w est de cardinal inférieur à b+1.

Il en résulte, dans le cas BI, en omettant l'indice supérieur ( $\infty$ ),

$$M_q(x) = p \cdot e^{((p-1)/p)x} \cdot M_q(x/p) + p \cdot (1/q!) (x/p)^q \cdot \sum_{j=b-q+1}^{\infty} (1/j!) \left(\frac{p-1}{p}x\right)^j$$

$$M_q(x) = e^x \cdot \sum_{j=0}^{\infty} p^j b_q(x/p^j) \cdot e^{-xp^{-j}} = p e^x \sum_{m=b+1}^{\infty} (x/p)^m \frac{(p-1)^{m-q}}{q!(m-q)!} \left( \sum_{k=0}^{\infty} p^{(1-m)j} e^{-xp^{-j}} \right)$$

Soit, en tenant compte du cas où w est de cardinal non supérieur à b,

$$M_q(x) = (x^q/q!) + p e^x \cdot \sum_{m=b+1}^{\infty} (x/p)^m \cdot \frac{(p-1)^{m-q}}{q!(m-q)!} \cdot S_m(x) .$$

Sous l'hypothèse P, on a  $M_q^{(P)}(v) = e^{-v} \cdot M_q(v)$ ; ce qui donne pour q=0,

$$M_0^{(P)}(v) = e^{-v} + p \cdot \sum_{m=b+1}^{\infty} (v/p)^m \frac{(p-1)^m}{m!} \cdot S_m(v) .$$

$$M_0^{(P)}(v) = e^{-v} + p \cdot \sum_{m=b+1}^{\infty} (v/p)^m \frac{(p-1)^m}{m!} \cdot v^{1-m} \frac{\Gamma(m-1)}{\text{Ln}(p)} [1 + N_m(\text{Log}_p v) + O(v^{-m})]$$

$$M_0^{(P)}(v) = ((pv)/\text{Ln}(p)) \cdot U_b \cdot [1 + W(\text{Log}_p v) + O(v^{-m})] , \text{ où } U_b = \sum_{m=b+1}^{\infty} \frac{1}{m(m-1)} u^m ,$$

$u = \frac{p-1}{p}$ , et W est une fonction périodique de  $\text{Log}_p v$ , de moyenne nulle et de période 1.

De plus 
$$U_b = \frac{1}{b(b+1)} u^{b+1} \left[ 1 + \sum_{m=1}^{\infty} \frac{b(b+1)}{(b+m)(b+m+1)} u^m \right]$$

Or  $u$  est positif et strictement inférieur à 1, et le coefficient de  $u^k$  est également inférieur à 1.

Par conséquent

$$\frac{1}{b(b+1)} u^{b+1} < U_b < \frac{1}{b(b+1)} u^{b+1} (1/(1-u)).$$

En résumé, on a le résultat suivant:

Théorème: asymptotiquement en  $v$ , sous l'hypothèse  $P$ , le nombre moyen de feuilles de l'index qui sont inoccupées,  $F_0^{(P)}(v)$ , est tel que

$$F_0^{(P)}(v) \sim \frac{pv}{\ln(p)} \cdot U_b, \text{ avec une précision de } O(v^{-m}), \text{ } m \text{ réel positif, et}$$

$$V_b < U_b < p \cdot V_b, \text{ en posant } V_b = \frac{1}{b(b+1)} \left(\frac{p-1}{p}\right)^{b+1}.$$

Remarque: du fait que  $p/\ln(p)$  et que  $\left(\frac{p-1}{p}\right)^{b+1}$  sont fonctions croissantes de  $p$ , on constate que le nombre de feuilles inoccupées de l'index est lui aussi fonction croissante de  $p$ ; ce qui était prévisible.

IV-3-Coût moyen d'insertion ou de suppression .

Nous avons établi que  $I^{(P)}(v) = (v^b/b!) \sum_{k=0}^{\infty} (1/p^k)^b \cdot e^{-v/p^k}$  .

Nous pouvons réécrire  $I^{(P)}(v) = (v^b/b!) \sum_{k=0}^{\infty} p^{(1-(b+1))k} \cdot e^{-v/p^k}$  .

Soit  $I^{(P)}(v) = (v^b/b!) \cdot S_{b+1}(v)$  .

En utilisant le lemme ,on obtient le

Théorème :

Asymptotiquement en  $v$  , le coût moyen d'insertion ou de suppression d'une clé dans un fichier est , sous l'hypothèse P ,

$$I^{(P)}(v) = \frac{v}{b \cdot \ln(p)} [ 1 + N_{b+1}(\text{Log}_p v) + O(v^{-m}) ] .$$

(Avec les mêmes notations que précédemment) .

On voit donc que ce coût est une fonction périodique de  $\text{Log}_p v$ , de période 1, et de valeur moyenne  $v/(b \cdot \ln(p))$ .

IV-4-Longueur du cheminement externe .

Désignons par  $L_b^{(P)}(v)$  ce nombre pour des arbres dont les boîtes associées aux feuilles sont de capacité  $b$ . Nous avons montré que

$$L_b(x) = e^x \sum_k p^k (x/p^k) \cdot (e^{x/p^k} - e_{b-1}(x/p^k)) \cdot e^{-x/p^k} . \text{ Soit}$$

$$L_b(x) = e^x \cdot x \cdot \sum_k (1 - e^{-x/p^k}) \cdot \left( \sum_{j=1}^{b-1} \frac{1}{j!} x^j \cdot p^{-jk} \right) \cdot e^{-x/p^k} \cdot x e^x . \text{ D'où}$$

$$L_b(x) = L_1(x) - x \cdot e^x \cdot \sum_{j=1}^{b-1} (x^j/j!) \cdot \sum_k p^{(0-j)k} \cdot e^{-x/p^k}$$

$$= L_1(x) - x \cdot e^x \cdot \sum_{j=1}^{b-1} (x^j/j!) \cdot S_{j+1}(x) \quad , \quad \text{avec}$$

$$L_1(x) = x \cdot e^x \cdot \sum_k (1 - e^{-x/p^k}) = x \cdot e^x \cdot T(x), \text{ avec } T(x) = \sum_k (1 - e^{-x/p^k}) .$$

Il résulte de ces égalités que

$$L_1^{(P)}(v) = v \cdot T(v), \quad \text{et que}$$

$$L_b^{(P)}(v) = L_1^{(P)}(v) - v \cdot \sum_{j=1}^{b-1} (v^j/j!) \cdot S_{j+1}(v) \quad , \quad \text{si } b \text{ est supérieur à } 1.$$

Du fait que l'on sait calculer les  $S_j(v)$ , la détermination des  $L_b^{(P)}(v)$  se ramène à celle de  $T(v)$ .

Lemme 2:

Asymptotiquement en  $x$ , on a

$$T(x) = (\text{Log}_p x + \frac{1}{2} + \frac{\gamma}{\text{Ln}(p)}) [ 1 + U(\text{Log}_p x) + O(x^{-m}) ] \quad ,$$

où  $U$  est une fonction périodique de  $\text{Log}_p x$ , de période 1, de moyenne nulle, et  $m$  un réel positif quelconque.

Preuve: elle est identique à celle de [1], en remplaçant 2 par  $p$ .

L'utilisation des deux lemmes conduit au

Théorème

Asymptotiquement en  $v$ , sous l'hypothèse  $P$ , la longueur du cheminement externe de l'index est une fonction périodique de  $\text{Log}_p v$ , de période 1, de valeur moyenne

$$L_b^{(P)}(v) \sim v [ \text{Log}_p v + 1/2 + (\gamma/\text{Ln}(p)) (1 + \delta_{p>1} \cdot \sum_{j=1}^{b-1} (1/j!)) ] .$$

ANNEXE

A titre d'exemple, étudions la série  $T(x) = \sum_{k=0}^{\infty} (1 - e^{-x/p^k})$ .

1) Sa convergence est assurée par le fait qu'asymptotiquement en  $k$

$$1 - e^{-x/p^k} \sim x/p^k, \text{ terme général d'une série convergente.}$$

2) Nous allons effectuer une transformation de Mellin sur  $T(x)$  :

$$\text{posons } \alpha_k = 1, \beta_k = p^{-k}, f(x) = 1 - e^{-x}.$$

$$\text{On a } T(x) = \sum_{k=0}^{\infty} \alpha_k \cdot f(\beta_k x).$$

Pour tout nombre complexe  $s$ , la transformée de Mellin  $T^{\times}(s)$  de  $T(x)$  est

$$T^{\times}(s) = \int_0^{\infty} T(x) \cdot x^{s-1} dx. \text{ Elle vérifie}$$

$$T^{\times}(s) = f^{\times}(s) \cdot \omega(s), \quad \omega(s) = \sum_k \alpha_k \cdot \beta_k^{-s} = \sum_k (p^{-k})^{-s} = \sum_k (p^s)^k, \text{ et}$$

$$f^{\times}(s) = -\Gamma(s).$$

$\omega$  est convergente pour  $\text{Re}(s) < 0$ , et vaut alors  $\omega(s) = 1/(1 - p^s)$ .

$f^{\times}(s)$  est définie, holomorphe pour  $-1 < \text{Re}(s) < 0$ , et y vaut  $-\Gamma(s)$ .

On a donc  $T^{\times}(s) = -\Gamma(s)/(1 - p^s)$ , pour  $-1 < \text{Re}(s) < 0$ .

Cette fonction peut être prolongée analytiquement pour  $\text{Re}(s) > 0$ .

3) pour tout réel  $c$ , compris entre  $-1$  et  $0$ , on a par inversion :

$$T(x) = (1/2i\pi) \int_{c-i\infty}^{c+i\infty} T^x(s) \cdot x^{-s} ds .$$

Prenons un contour C pour évaluer cette intégrale: soit m un réel positif, et  $t_n = (2\pi+1)n/\ln(p)$ , n entier positif tendant vers l'infini. Le contour est délimité par les droites d'abscisses c et m, et les droites d'ordonnées  $t_n$  et  $-t_n$ . Il est parcouru dans le sens direct. On désigne par  $I_c$ ,  $I_m$ ,  $I_{sup}$  et  $I_{inf}$  les intégrales de  $T^x(s) \cdot x^{-s}$  sur ces différentes côtés de ce rectangle parcourus dans le sens direct, et par B l'intérieur de ce rectangle. La fonction à intégrer est analytique à l'intérieur du contour, sauf en ses pôles que nous étudierons plus loin. On peut donc lui appliquer le théorème de Cauchy :

$$I_c + I_m + I_{sup} + I_{inf} = \sum_{s_i \in S} \text{Res}(T^x(s) \cdot x^{-s}) ,$$

où S est l'ensemble des pôles  $s_i$  de  $T^x(s) \cdot x^{-s}$  qui se trouvent dans B.

On constate que  $T(x) = -\lim I_c$ , pour n tendant vers l'infini.

Nous allons nous employer à calculer la somme des résidus et à montrer le rôle négligeable des autres intégrales.

4) Pôles et résidus de  $T^x(s) \cdot x^{-s}$  dans B:

$x^{-s}$  n'ayant pas de pôle dans B, les seuls pôles à prendre en considération sont ceux de  $1/(1-p^s)$  et ceux de  $\Gamma(s)$ .

Les pôles de la première de ces fonctions vérifient

$$p^s = 1, \text{ soit } e^{s \cdot \ln(p)} = e^{2il\pi}, \text{ l entier relatif. Donc ce sont les}$$

$$s_l = 2il\pi/\ln(p), \text{ } l \in \mathbb{Z}.$$

A l'intérieur de B,  $\Gamma(s)$  n'a qu'un seul pôle,  $s = 0$ .

La fonction  $T^x(s) \cdot x^{-s}$  a donc un pôle double à l'origine et une infinité de pôles simples  $s_1 = 2i\pi/\ln(p)$ ,  $l \in \mathbb{Z}^*$ .

Résidu à l'origine:

au voisinage de  $s=0$ , on a  $\Gamma(s) = \Gamma(s+1)/s \sim (1 - \gamma \cdot s)/s$ , où  $\gamma$  est la constante d'Euler, ([4], p.257 et suivantes).

$$1 - p^s = 1 - e^{s \cdot \ln(p)} \sim -s \ln(p) \left(1 + \frac{s}{2} \ln(p)\right);$$

$$x^{-s} = e^{-s \cdot \ln(x)} \sim 1 - s \cdot \ln(x).$$

$$\text{Donc } T^x(s) \cdot x^{-s} \sim (1/s^2 \ln(p)) \cdot \left(1 - \left(\gamma + \frac{1}{2} \ln(p) + \ln(x)\right)s\right).$$

Le résidu en  $s=0$  est le coefficient de  $1/s$ , soit

$$\rho_0 = -\left(L_p(x) + 1/2 + \gamma/\ln(p)\right).$$

Résidu en chacun des autres pôles  $s_1$ :

$$\rho_1 = -\Gamma(s_1) \cdot x^{-s_1} / (1 - p^{s_1})'_{s=s_1} = \frac{\Gamma(s_1)}{\ln(p)} \cdot x^{-s_1}.$$

5) Etude de  $I_{\text{sup}}$  et  $I_{\text{inf}}$ :

Sur la droite d'équation  $s = i \cdot t_n$ , on a,

$$1 - p^s = 1 - p^{s_n} \cdot e^{in/\ln(p)} = 1 - e^{in/\ln(p)} = 1 - \cos(n/\ln(p)) + i \sin(n/\ln(p))$$

Donc  $|1 - p^s| = 2 \cdot \sin(n/2 \cdot \ln(p))$ , qui ne peut nul pour des raisons

évidentes. Il existe donc un réel positif  $d$  tel que  $d \leq |1 - p^s| \leq 2$ .

De plus, on sait que, quand  $v$  tend vers l'infini, on a

$$\Gamma(u+iv) = O(|v|^{u-1/2} \cdot e^{-(\pi/2)|v|}) .$$

$$\text{Donc } \Gamma(it_n) = O(|t_n|^{1/2} \cdot e^{(\pi/2)|t_n|}) .$$

$$\text{Enfin } x^{-s} = x^{-t_n} = e^{-it_n \text{Ln}(x)} \text{ et } |x^{-s}| = 1 .$$

Par conséquent  $T^x(s) \cdot x^{-s}$  est majorée en module par une fonction de  $t_n$  qui tend vers 0 quand  $n$  tend vers l'infini, et l'on a

$$\lim I_{\text{sup}} = \lim I_{\text{inf}} = 0 , \text{ quand } n \text{ tend vers l'infini.}$$

6) Etude de  $I_m$  :

c'est l'intégrale de  $T^x(s) \cdot x^{-s}$  sur le côté  $s = m + iy$ ,  $y$  allant de  $-t_n$  à  $+t_n$ . On a, ici encore,

$$|\Gamma(m+iy)| = O(|y|^{m-1/2} \cdot e^{-(\pi/2)|y|}) ;$$

$$p^s = p^m \cdot e^{iy \text{Ln}(p)} , \quad 1 - p^s = 1 - p^m \cdot \cos(y \text{Ln}(p)) - i p^m \sin(y \text{Ln}(p)) , \text{ donc}$$

$$p^m - 1 \leq |1 - p^s| \leq p^m + 1 ;$$

$$x^{-s} = x^{-m} \cdot e^{-iy \text{Ln}(x)} , \text{ soit } |x^{-s}| = x^{-m} , \text{ pour } x \text{ positif .}$$

Il s'ensuit que

$$|I_m| \leq (x^{-m}/(p^m-1)) \cdot |\int \Gamma(m+iy)| .$$

Or cette dernière intégrale converge vers une limite finie , quand n tend vers l'infini, compte tenu de la majoration obtenue pour  $|\Gamma(m+iy)|$  :

Lorsque x tend vers plus l'infini, on voit donc que le module de  $I_m$  est un  $O(x^{-m})$  , quelque soit le réel positif m.

On a donc, en résumé, par passage à la limite, quand n tend vers l'infini,

$$T(x) = -\rho_0 + \sum_{1 \in \mathbb{Z}^x} \rho_1 + O(x^{-m}) , \text{ m réel positif quelconque, soit encore}$$

$$\sum_{k=0}^{\infty} (1 - e^{-x/p^k}) = L_p(x) + 1/2 + \gamma/\ln(p) + \sum_{1 \in \mathbb{Z}^x} (1/\ln(p)) \Gamma(2i1\pi/\ln(p)) e^{-2i1\pi \text{Log}_p x} + O(x^{-m}) ; \text{ ou enfin}$$

$$\sum_{k=0}^{\infty} (1 - e^{-x/p^k}) = \text{Log}_p x + 1/2 + \gamma/\ln(p) + V(\text{Log}_p x) + O(x^{-m}),$$

où V est une fonction périodique de  $\text{Log}_p x$  , de période 1 et de moyenne nulle , et m un réel positif quelconque .

## CONCLUSION

Ainsi que nous l'avions indiqué, nous avons généralisé au cas des arbres  $p$ -homogènes les résultats connus pour le hachage dynamique dans le cas des arbres binaires.

deux constatations s'imposent:

d'une part, méthodologiquement, les outils créés et utilisés pour  $p=2$  ont pu être généralisés sans trop de difficulté à  $p$  entier quelconque; d'autre part, au plan des résultats, le fait de prendre  $p$  supérieur à 2 n'améliore pas les performances du hachage. En effet, comme on l'a vu, les valeurs de certains des paramètres se dégradent quand  $p$  augmente. Néanmoins, sauf à se résigner à coder en binaire les clés écrites dans un alphabet à  $p$  éléments, on peut considérer que les procédés de hachage au moyen d'arbres  $p$ -homogènes sont très utiles.

Pour terminer, nous voudrions suggérer une orientation de recherche, qui, à notre avis, devrait apporter une amélioration des performances des techniques de hachage dynamique : alors que les méthodes actuelles consistent à séparer, de façon aveugle, les clés qui sont rangées dans une boîte attachée à une feuille de niveau  $k$  dans l'index, suivant les valeurs de leurs  $k$ -ièmes chiffres, on pourrait penser à les séparer suivant les valeurs d'un chiffre - pas nécessairement le  $k$ -ième - qui partagerait "le mieux possible", en un sens à préciser, cet ensemble de clés.

## BIBLIOGRAPHIE

- [1] M. Régnier, "Evaluation des performances du hachage dynamique", Thèse, Université Paris-Sud (1983).
- [2] P. Flajolet, M. Régnier, R. Sedgewick : "Some uses of the Mellin integral transform in the analysis of algorithms". Proc. Advanced Workshop on "Combinatorics on Words", Springer NATO ASI Ser. F12, pp 241-254.
- [3] P. Flajolet, M. Régnier, D. Sotteau : "Algebraic methods for trie statistics in analysis and design of algorithms for combinatorial problems". Annals of Discrete Math., Vol 25-1985 pp 145-188.
- [4] Abramovitz and Stegun : "Handbook of mathematical functions".
- [5] R. Fagin., J. Nievergelt, N. Pippenger., H.R. Strong. "Extendible Hashing - a fast access method for dynamifiles". ACM-TODS, 4, 3, (Sep. 1979), 315-344.
- [6] P.A. Larson., "Dynamic Hashing" BIT 18, (1978), 184-201.
- [7] W. Litwin., "Virtual Hashing : a dynamically changing hashing". VLDB 78. ACM, (Sept. 1978), 517-523.

Imprimé en France

par

l'Institut National de Recherche en Informatique et en Automatique

