

MINISTERE DE L' EDUCATION NATIONALE

SECRETARIAT D' ETAT A LA RECHERCHE

**THESE**

Présentée au

CENTRE DE DEVELOPPEMENT DES TECHNOLOGIES AVANCEES

Pour l'obtention du grade de

**MAGISTER EN CYBERNETIQUE**

( Option : Architecture des systèmes )

**OPTIMISATION DU STOCKAGE ET DE L' ACCES  
DANS LES BASES DE DONNEES  
TRES VOLUMINEUSES**

Par

**Mme Fatma-Zohra BESSAI  
née MECHMACHE**

Soutenue le 45 Novanbre 1992 Devant le jury composé de

Mr H. Bessalah	.....	Président
Mr A. Abdellaoui	.....	Examineur
Mme Z. Ali-Mazighi	.....	Examineur
Mr M. Benhamadi	.....	Examineur
Mme F. Cherief	.....	Examineur
Melle F. Azrou	.....	Rapporteur

## Remerciements

Je tiens a remercier,

Monsieur H. BESSALAH, Directeur du centre de developpement des technologies avancees (CDTA), pour avoir bien voulu me faire l'honneur de presider le jury.

Monsieur M. BENHAMADI, Directeur du CERIST, pour avoir mis a ma disposition les moyens necessaires a la bonne marche de mon travail, pour la confiance qu'il m'a temoignee, et pour sa participation au jury.

Je suis tres honoree de la participation a ce jury de Madame Z. ALI-MAZIGHI, chargee de cours a l'USTHB, Madame F. CHERIEF, maitre de conference (USTHB) et de Monsieur A. ABDELLAOUI, maitre de conference (USTHB).

Je tiens egalement a remercier Mademoiselle F. AZROU, chargee de recherche au CERIST, pour avoir propose le sujet de recherche et dirige mon travail.

Je ne peux pas oublier ici H. LABIOD, Ingenieur au CERIST, pour son aide et sa collaboration pendant l'integration de mon travail au systeme SIBADOC.

Je ne saurais clore ces pages sans adresser mon expression de sympathie la plus sincere a tous mes collegues et amis du Laboratoire de Recherche et Developpement en Informatique (LRDI), ainsi qu'a l'ensemble du personnel du CERIST. Que toutes celles et ceux qui m'ont temoigne leur disponibilite et apporte leur aide trouvent ici une expression de reconnaissance.

Enfin je tiens a remercier tous mes proches pour leur soutien et patience tout au long de l'elaboration de ce travail. Qu'ils trouvent ici une recompense a leurs efforts.

## **RESUME**

La masse d'information mise a la disposition de la communauté scientifique s'accroît de façon continue, étant donnée l'évolution et la diversification des disciplines scientifiques et techniques. Cette grande masse d'information nécessite l'optimisation des moyens de stockage et le développement de nouvelles méthodes d'accès car les méthodes classiques se révèlent insuffisantes.

L'objectif de notre travail est de contribuer à l'amélioration des techniques d'accès aux bases de données volumineuses. Or une méthode d'accès est étroitement liée à la méthode de stockage utilisée. Aussi, une partie de notre travail a consisté à étudier le compactage de données, utilisé au CERIST pour résoudre le problème de stockage de grands volumes de données.

Après l'étude des diverses techniques de compression, notre travail a essentiellement consisté à concevoir et à réaliser un système d'accès aux données volumineuses. Ce système utilise une technique d'indexation basée sur l'inversion totale, ce qui permet aux utilisateurs d'accéder à la base de données par n'importe quel mot contenu dans cette dernière et offre un accès rapide à l'information grâce à l'utilisation du Q-arbre qui est un arbre équilibré de profondeur 1. En plus de ces caractéristiques, le système a une structure modulaire lui permettant de s'adapter à tous les systèmes gérant de grands volumes de données textuelles.

## **MOTS CLES**

Base de données documentaire, système de recherche documentaire, système de gestion d'objets complexes, stockage des données, accès aux bases de données, compression de données, indexation, Q-arbre.

## **ABSTRACT**

The amount of scientific information is growing steadily. This huge amount of information requires efficient storage management and new access methods.

Our goal is to improve the current access techniques to very large data bases. Since an access technique is closely related to the storage method used, in the first part of our work we study the data compression techniques already used by the CERIST in order to store a large amount of information.

After an analysis of the current techniques of data compression, we concentrate on the design and implementation of an access system to very large data bases. Our system uses an indexing scheme based on full inversion which allows users to query the data base using any word contained in the data base. This system is based on the Q-tree (B-tree of depth one) that indeed allow fast access to the sought information. In addition, the system modular structure makes it adaptable to all system managing very large textual data.

## **KEY WORDS**

Bibliographic data base, document-retrieval system, complex object management system, data storage, data base access, data compression, indexing, Q-tree.

## S O M M A I R E

### Introduction

I	Evolution des systemes de gestion de donnees	
1.1	Les systemes de recherche documentaire .....	5
1.2	Les systemes de gestion de bases de donnees .....	7
II	Compression de donnees	
11.1-	Definition .....	13
11.2-	Concepts de base de la compression de donnees (CD) .	13
11.2.1-	Le code et ses proprietes .....	14
11.2.2-	Quantité d'information et optimalité .....	16
	d'un code	
11.2.3-	Schema de compression de donnees .....	17
11.2.4-	Parametre d'évaluation des techniques .....	19
	de compression de donnees (TCD)	
11.3-	Les redondances dans la représentation des donnees .	19
11.4-	Classification des techniques de compression de ....	20
	donnees	
11.5-	Choix d'une TCD pour la compression des donnees ....	25
	bibliographiques	
11.6-	Integration de la technique de compression de .....	29
	donnees au systeme SIBADOC	
III	Techniques d'accès	
111.1	Description des differentes techniques d'accès ...	37
III.1.1	Recherche sequentielle .....	37
	(Full Text Scanning)	
III.1.2	Inversion .....	37
III.1.3	Regroupement (Clustering) .....	38
III.1.4	Fichier signature (Signature file) .....	39
III.1.5	Methode d'accès utilisant un fichier ....	40
	signature	
III.1.6	Methode d'accès utilisant un Q-arbre ....	43
III.1.6.1	Definition d'un Q-arbre .....	43

	III.1.6.2	Description du Q-arbre .....	44
	111.2	Criteres de choix de la technique d'accès .....	45
IV		Indexation ou selection de mots cles	
	IV.1	Les methodes d'indexation .....	48
		IV.1.1	Indexation sans thesaurus et indexation .... 49
			avec thesaurus
		IV.1.1.1	Definition d'un langage .....
			documentaire
		IV.1.1.2	Definition d'un thesaurus .....
		IV.1.1.3	Indexation sans thesaurus .....
		IV.1.1.4	Indexation avec thesaurus .....
		IV.1.2	Indexation manuelle et indexation .....
			automatique
		IV.1.2.1	Indexation manuelle .....
		IV.1.2.2	Indexation automatique .....
			IV.2.2.2.1
			Approche statistique . 54
			IV.2.2.2.2
			Approche linguistique 55
	IV.2	Criteres de choix de la methode d'indexation .....	56
		IV.2.1	Les systemes en texte integral ou .....
			systemes en "full text"
		IV.2.2	Methodes d'indexation automatique .....
			IV.2.2.1
			Methode statistique .....
			IV.2.2.2
			Methode linguistique .....
	IV.3	Conclusion .....	60
V		Systeme d'accès aux donnees volumineuses	
	V.1	Description du systeme d'accès aux donnees .....	63
			volumineuses
		V.1.1	Module d'indexation .....
		V.1.2	Module d'accès aux donnees .....
		V.1.3	Module de mise a jour de la Base de .....
			Donnees Documentaires (BDD)
	V.2	Integration du systeme d'accès au systeme SIBADOC ....	68
		V.2.1	Module d'indexation .....
			V.2.1.1
			Strusture du fichier lexique .....
			V.2.1.2
			Structure du fichier inverse .....

v.2.2	Module d'accès aux données .....	71
v.2.3	Module de mise à jour .....	73
v.3	Evaluation du système d'accès aux données ..... volumineuses	73
v.3.1	Evaluation de l'indexation .....	73
v.3.2	Evaluation de la méthode d'accès .....	77
	Conclusion et perspectives .....	80

## Bibliographie

Annexe A: SIBADOC : Un système d'interrogation des bases de données documentaires

Annexe B: Liste des mots vides utilisés par le module d'indexation du système d'accès aux données volumineuses.

.

## ***INTRODUCTION***

Depuis l'apparition de la notion de bases de données vers le début des années 60, on assiste à une prolifération de systèmes de gestion de données opérant sur de grandes quantités d'information. Les possibilités limitées de stockage et d'accès rendent difficile la manipulation de cette importante masse de données. Par conséquent, à défaut de disposer de toutes les capacités de stockage nécessaires, seules les données les plus récentes sont accessibles aux utilisateurs. Ce problème est posé de façon aiguë, au Centre de Recherche sur l'Information Scientifique et Technique (CERIST), dans les deux contextes suivants:

- D'une part, lors de l'exploitation des bases documentaires internationales telles que INSPEC, INIS et AGRIS. En effet, à titre d'exemple, une base annuelle de INSPEC nécessite environ 300 MegaOctets de mémoire, sans compter les informations de gestion utilisées par le SRD. Aussi actuellement, seule la base qui correspond à la dernière année sera accessible aux utilisateurs.

- D'autre part, lors du développement du système de gestion d'objets complexes SYGOC [BOU 91]. En effet, le système SYGOC est conçu pour la gestion et le stockage des données volumineuses et à structure complexe, tel que les documents, d'où nécessite de beaucoup d'espace et de méthodes d'accès optimisées.

Ce problème de stockage et de manipulation de grands volumes de données nous a incité à développer un nouveau système informatique, dont le but est de manipuler efficacement des bases de données très volumineuses, et ceci en :

- Minimisant l'espace de stockage pour les bases de donnees.
- Permettant une recherche rapide.
- Effectuant les mises a jour des bases de donnees en un temps satisfaisant.

Pour reduire l'espace de stockage de la base, une solution a ete adoptee, c'est la compression de donnees. Ce travail a été déjà réalisé et a conduit au developpement d'un schema de compression de donnees [MEC 89], qui, intégré au systeme de recherche documentaire **SIBADOC** [BEN 88], developpe au CERIST, a donne des resultats satisfaisants. En effet, la base INSPEC est réduite a environ le tiers de sa taille initiale etant donne que le taux de compression est approximativement égal a 2/3.

Le probleme de l'optimisation dans l'espace de stockage des donnees ayant ete ainsi resolu, notre travail consistera a developper un systeme d'accès rapide qui pourrait acclereler la recherche de l'information, et par consequent, augmenter l'efficacité de notre systeme de stockage et d'accès dans des bases de donnees volumineuses.

En ce qui concerne la realisation, le systeme SYGOC etant en cours de developpement, nous nous sommes interessees a l'application au systeme de recherche documentaire **SIBADOC** qui lui est en phase d'exploitation. Par la suite tout ce qui a été developpe sur **SIBADOC**, pourrait être transpose sur le systeme **SYGOC**, moyennant quelques modifications, car les problemes **poses** par les deux systemes sont identiques.

Après une étude de l'évolution des systèmes de gestion de données, présentée dans le 1er chapitre, le reste du document est organisé de la manière suivante:

Le 2ème chapitre présente une étude sur la compression de données, en donnant une double classification des techniques de compression de données (TCD), tout en présentant les TCD les plus couramment utilisées et les critères de choix d'une TCD pour la compression de données bibliographiques, ainsi que l'intégration de cette dernière au système SIBADOC.

Dans le 3ème chapitre nous mettons en évidence les différentes méthodes d'accès aux données ainsi que les principaux critères qui permettent de comparer ces méthodes entre elles. Nous présentons ensuite, les critères de choix de la technique d'accès utilisée dans notre système, pour rechercher l'information.

L'efficacité des méthodes d'accès aux données dépend en général de l'indexation de ces données. C'est-à-dire du choix des termes permettant de les retrouver. Aussi, le 4ème chapitre est-il consacré à l'étude des différentes méthodes d'indexation, tout en présentant les critères de choix de la technique d'indexation retenue pour notre système.

Le 5ème chapitre est dédié aux réalisations et expérimentations des techniques d'accès et d'indexation que nous avons définies. Nous y présentons le système d'accès aux données volumineuses, ainsi que l'intégration de ce dernier au système SIBADOC.

# ***CHAPITRE I***

***Evolution des systèmes de gestion de données***

Afin de situer le contexte dans lequel notre travail a été réalisé, nous abordons dans ce chapitre les thèmes suivants :

- Les Systèmes de Recherche Documentaire (SRD)
- Les Systèmes de Gestion de Base de Données (SGBD).

### 1.1- Les systèmes de recherche documentaire

La recherche documentaire englobe l'ensemble des techniques et modalités permettant de sélectionner l'information dans un fonds documentaire moyennant des critères spécifiés par l'utilisateur [DEW 83] [DEW 89]. Son automatisation a été entreprise pour différentes raisons parmi lesquelles nous citons :

- Impossibilité de canaliser le flux d'information à partir d'un certain seuil.
- Impossibilité de maîtriser les volumes de données des fonds documentaires, s'accroissant linéairement avec le temps.
- Impossibilité d'assurer une certaine précision informative à la recherche.
- Impossibilité d'accéder rapidement à l'information **des** lors que le volume couvrant l'information recherchée devient important et que la question devient complexe.

Cette informatisation de la recherche documentaire a conduit au développement de systèmes d'information spécialisés dans le stockage et la manipulation de données bibliographiques : Ce sont les systèmes de recherche documentaire.

La fonction essentielle des systemes de recherche documentaire consiste à stocker les documents sous forme d'un ensemble d'information caractérisant ces derniers, et à en permettre l'accès lors des recherches bibliographiques effectuées par les utilisateurs du systeme.

Les informations manipulees par le systeme ne sont pas les documents eux-mêmes, mais une description de ces derniers. Generalement, les systemes de recherche documentaire representent un document à l'aide d'une liste bibliographique, apelée notice bibliographique, et qui regroupe des informations telles que: nom d'auteur, references du document, titre, resume, annee d'édition, lieu d'édition, etc.

Les **SRD** gerent de grands volumes de donnees, ayant une representation condensee, comme par exemple des resumes associes à des listes de mots-cles appartenant à des vocabulaires contrôlés. Dans un **SRD** les insertions de documents sont realises par "**lots**", par un administrateur du systeme et les suppressions sont tres rares, pour ne pas dire inexistantes. Les mises à jour de documents ne se presentent jamais.

L'offre en **SRD** connaît un veritable essor. Chaque annee des dizaines de systemes sont developpes. Parmi ces systemes nous avons le systeme **STAIRS** [IBM 79], **DIAMG** [BOU 79], **SMART** [SAL 71] [SAL 86a], etc. Helas, malgre la croissance rapide du nombre de ces systemes, la plupart d'entre eux, pour ne pas dire tous, utilisent la technique d'inversion classique pour acceder aux donnees bibliographiques.

## 1.2- Les systemes de gestion de base de donnees

Un systeme de gestion de base de donnees est aujourd'hui un logiciel de base essentiel dans un systeme informatique de gestion. Intuitivement, il permet a des utilisateurs concurrents de manipuler (insérer, modifier et rechercher) efficacement des donnees contenues dans une base de donnees. Depuis l'apparition des premiers SGBD vers 1962, d'importants resultats theoriques et pratiques ont ponctue l'histoire de la recherche en base de donnees. La mise en oeuvre de ces resultats a permis de faciliter l'administration et la manipulation d'une base de donnees, et d'accroître ainsi la productivite des utilisateurs des bases de donnees.

Les SGBD sont generalement classes selon le modele de donnees sur lequel ils sont bâtis. Leur histoire peut etre résumée en distinguant trois generations.

La premiere generation des SGBD est basee sur les modeles RESEAU et HIERARCHIQUE où les langages de manipulation des donnees sont de type procedural (navigationnel), autrement dit, l'utilisateur doit specifier le chemin d'accès aux donnees afin d'obtenir l'information desiree.

L'idée de la deuxième generation de SGBD est nee vers 1970 avec l'apparition du modele de donnees relationnel, défini par E.F CODD et base sur des fondements theoriques solides [COD 70] [GAR 85] [GAR 91]. Ceci a ouvert le champ a de nombreuses recherches de par le monde en vue de construire un systeme repondant aux normes du modele relationnel. Parmi ces projets deux systemes peuvent être consideres comme des references, il s'agit de SYSTEM-R developpe dans les laboratoires d'IBM [AST 76] et de INGRES chez BELL laboratory [STO 76].

Cependant, les premiers systemes relationnels n'ont ete commercialises que vers le debut des annees 80, avec l'apparition sur le marche du systeme SQL/DS [GAR 85] issu de **SYSTEM-R**.

Le modele relationnel temoigne d'une evolution des SGBD en parvenant a une plus grande independance des donnees par rapport aux programmes. Cette independance est d'une part physique : les applications ne doivent pas dependre du mode de stockage des donnees, et d'autre part logique: les utilisateurs doivent pouvoir conserver leur vision propre des donnees contenues dans la base.

Les SGBD relationnels sont tres bien adaptes pour les applications de gestion classique, comme , par exemple, la comptabilite, la gestion de stock ou le controle d'inventaire. Cependant une tendance recente souligne les limites de ces systemes. Cette tendance est la manifestation de besoins pressants en gestion de donnees de la part d'applications nouvelles, telles que la bureautique, la conception assistee par ordinateur (**CAO**), l'exploitation des bases de donnees specifiques ou le genie logiciel, dont le besoin commun est la manipulation d'objets complexes (textes, graphiques, cartes, images, donnees multidimensionnelles), c'est-a-dire, des objets dont les attributs ne sont pas necessairement atomiques mais peuvent etre eux memes des objets. Autrement dit, des donnees de tres grande taille ayant une semantique riche et une structure complexe.

Par exemple, soit l'objet complexe LIVRE constitue d'un ensemble de chapitres et chaque chapitre contenant un ensemble de paragraphes. De tels objets peuvent etre aisement modelises avec le modele relationnel en regroupant les objets de meme type dans une meme relation et en associant les sous-objets d'un meme objet par des attributs de jointure.

La modelisation de l'objet complexe LIVRE nécessite alors, de façon tres simplifiée, les relations suivantes :

```
LIVRE(cote,titre,auteur) ;
CHAPITRE(num-chap,cote,titre-chap) ;
PARAGRAPHE(num-parag,num-chap,...)
```

ou cote et num-chap, respectivement, dans les relations CHAPITRE et **PARAGRAPHE** sont des attributs de jointure.

Un premier probleme avec l'approche relationnelle est que les liens sémantiques entre sous-objets sont perdus et doivent être explicites par des contraintes d'intégrité dont le maintien reste difficile et exige des algorithmes sophistiqués [VAL 87].

Nous devons specifier dans le cas de notre exemple, que "cote" dans la relation **CHAPITRE** est une cle étrangere, ce qui garantit qu'un chapitre n'existe pas s'il n'est pas rattache a un livre particulier.

Un autre problème qui se pose, est que les performances de la manipulation des objets complexes sont faibles et souvent inacceptables. Par exemple, la reconstitution des informations concernant un livre exige, dans notre cas, deux jointures (operation coûteuse) des trois relations precitees. Generalement le nombre de jointures est tres important d'où un temps d'accès prohibitif.

Le support d'objets complexes par les **SGBD** relationnels actuels est difficile pour les raisons suivantes:

1- Les types de donnees supportes sont en general limités a quelques domaines numériques ou alphanumériques ( par exemple, entier, reel et chaîne de caracteres).

**2- Perte de sémantique des objets complexes:** Les structures de données fournies par le modèle relationnel sont trop simples pour modéliser des objets à structure complexe comme une hiérarchie ou un graphe. Un objet complexe (par exemple, un objet hiérarchique de CAO) peut être facilement décomposé en relations. Cependant, cette solution pose un problème sérieux. L'information sémantique véhiculée par la structure complexe d'un objet est dispersée sous forme de valeurs dans différentes relations, imposant à l'utilisateur d'adapter sa vision à celle du SGBD.

**3- Dégradation des performances lors de la manipulation des objets complexes:** la composition de données d'un même objet complexe peut exiger un grand nombre de jointures (une opération coûteuse en temps d'exécution).

Le modèle relationnel étant incapable de prendre en compte la sémantique des objets complexes, les chercheurs se sont tournés vers une troisième génération de SGBD, qualifiés aussi d'avancés, qui supportent bien ces applications nouvelles ou exploitent des environnements opérationnels complexes. Cette nouvelle génération, initiée dans quelques laboratoires de recherche vers la fin des années 70, fait maintenant l'objet de travaux de recherche et de développement intenses.

Les SGBD avancés ont pour ambition de lever les limites des SGBD relationnels actuels. Ils suscitent aujourd'hui de nombreux travaux de recherche et de développement. Cependant, deux approches importantes peuvent être isolées: les systèmes de gestion de bases de données orientées objets et les systèmes de gestion de bases de données relationnelles étendues.

- Les SGBD orientes objets [BAN 88] [GAR 91] ont pour objectif essentiel le support d'objets complexes. Cette approche consiste a intégrer les techniques des langages de programmation orientes objets (comme C++) et des bases de donnees afin d'offrir un langage unique et general pour programmer les applications traditionnelles et recentes des bases de donnees. De nombreux systemes orientes objets sont déjà operationnels. Parmi ces systemes nous avons les systemes O2 [LEC 88], ORION [KIM 87] et GBASE [GAR 91].

- Les SGBD relationnels etendus [GAR 89] [COD 79] cherchent a atteindre les objectifs des bases de donnees orientées objets tout en conservant la puissance et le support theorique du modele relationnel.

Dans le cadre des systemes gerant ces nouvelles bases de donnees, et des systemes de recherche documentaire, deux aspects paraissent essentiels: l'optimisation de l'espace de stockage et celle du temps d'accès a l'information qui font l'objet de notre travail.

## ***CHAPITRE II***

### ***COMPRESSION DE DONNEES***

### **11.1- Definition**

On désigne par compactage de données toute technique qui réduit la taille de la représentation physique des données tout en préservant un sous-ensemble d'informations. Par définition, une technique de compactage est dépendante de la sémantique des données.

La compression de données est une technique de compactage complètement réversible et qui a pour objectif la minimisation des quantités de données à stocker [GOT 75]. Elle est réalisable par un changement de représentation (codage) des données.

La compression de données offre de nombreux avantages parmi lesquels nous pouvons citer:

- La réduction des coûts de stockage et de transmission.
- L'amélioration des performances des systèmes de gestion de l'information (réduction du nombre d'entrées/sorties physiques).

En contrepartie, l'utilisation de la compression de données implique certains inconvénients d'implémentation propres à des techniques particulières telles que, la manipulation de bits, la représentation à longueur variable des données, la maintenance du code, etc.

### **11.2- Concepts de base de la compression de données**

La compression de données est une branche de la théorie de l'information qui a pour objectif la minimisation des quantités de données à transmettre. Celle-ci est réalisée par un changement de représentation (codage) des données.

Les elements de base de la theorie de l'information et de la codification sont presentes dans ce qui suit.

### II.2.1- Le code et ses proprietes

Un code est une application d'un ensemble de messages sources, écrits dans un alphabet source, vers un ensemble de mots de code (mots-code) écrits dans un alphabet code [HEA 72].

**Exemple:** Le code ASCII

Alphabet source = {a..z, A..Z, 0..9, ...}

Alphabet code = {0,1}

Ensemble de mots de code = (0000000 .. 1111111)

Un code est entièrement défini par la donnée de trois elements;

- L'ensemble des messages sources a coder,
- L'ensemble des mots-code,
- La fonction de correspondance entre les elements de ces deux ensembles.

Les codes sont classes suivant la longueur des messages sources est des mots-code. On en distingue ainsi quatre categories:

- Bloc-bloc
- Bloc-variable
- Variable-bloc
- Variable-variable

- Bloc: Le message source (resp. mots-code) a une longueur fixe, la même pour tous les éléments. Elle peut être égale à 1 caractère, 2 caractères ou plus (resp. un octet, 12 bits).

- Variable: La longueur des messages sources (resp. mots-code) est variable. Elle peut être une chaîne quelconque de caractères (fragment, mot) (resp. une chaîne de bits, de quartets ou d'octets).

### - Propriétés d'un bon code

1- Un code est distinct si tous les mots de code peuvent être distingués.

2- Un code est décodable d'une façon unique si dans toute séquence de mots de code le décodage d'un élément ne donne lieu à aucune ambiguïté.

3- Un code a la propriété préfixe si aucun mot-code n'est préfixe d'un autre. Sa décodabilité unique est assurée sans recourir à des informations supplémentaires (autres que le mot-code en cours de décodage) pour lever les ambiguïtés dans le décodage d'une chaîne.

4- Un code préfixe est minimal si pour tout  $x$ , si  $x$  est un préfixe d'un mot-code alors pour tout  $y$  appartenant à l'alphabet code,  $xy$  est un mot-code ou préfixe d'un mot-code. Cette propriété implique que le code a une longueur minimale.

**Exemples :**

1- Le code  $C: A \longrightarrow B$  ( $A=\{a,b,c\}$ ;  $B=\{1,10000,100\}$ ) est decodable de façon unique, mais le decodage d'une sous-chaîne commençant par 1 necessite la connaissance de ce qui suit le code pour lever l'ambiguité dans le decodage.

2- Le code  $C: A \longrightarrow B$  ( $A=\{a,b,c\}$ ;  $B=\{1111,1010,00\}$ ) verifie la propriete prefixe, aucune ambiguite ne peut apparaitre dans le decodage,

3- Le code  $C: A \longrightarrow B$  ( $A=\{a,b,c\}$ ;  $B=\{00,10,01\}$ ) verifie la propriete prefixe mais n'est pas minimal car 1 est prefixe de 10 alors que 11 n'est ni mot-code ni prefixe d'un mot-code. Le code minimal serait  $\{00,01,1\}$ .

**11.2.2- Quantite d'information et optimalité d'un code**

La quantité d'information apportee par la réalisation d'un evenement  $a_i$  (dans notre cas,  $a_i$  est l'occurrence d'un message dans le texte a comprimer) suivant une probabilité d'occurrence  $p(a_i)$ , est mesuree par:  $-\text{Log}_2 p(a_i)$  [HEA 72].

Le contenu moyen en information d'un message d'une source (texte a comprimer) est donne par la sommation des quantités d'information apportees par chaque message ponderees par leurs probabilites d'occurrence:

$$H = - \text{SOM}(\text{pour } i=1 \text{ a } N) (p(a_i) \text{Log}_2(p(a_i))) \quad [\text{HEA } 721]$$

ou  $N =$  nombre de messages distincts dans la source.

$H$  est appelee l'entropie de la source.

L'optimalité d'un code est considérée dans le sens où sa redondance est minimale. La redondance d'un code est définie par la différence entre sa longueur moyenne ( $L$ ) et l'entropie de la source qu'il codifie [HEA 72].

$$R = L - H$$

$$\text{où } L = \text{SOM}(\text{pour } i=1 \text{ à } N) (p(a_i) * l_i)$$

$$l_i = \text{longueur du code associé au message } a_i.$$

L'efficacité d'un code est définie comme le complément de la redondance par rapport à  $L$ ,

$$E = H/L$$

Un code est redondant si on a  $R > 0$ , il est dit à redondance minimale si  $R$  est minimale (la longueur moyenne du code ( $L$ ) proche de l'entropie ( $H$ ))

**Remarque :**

Pour faire tendre  $L$  vers  $H$  (et donc minimiser  $R$ ) il est préférable d'assigner des codes courts aux éléments les plus fréquents et de longs codes aux éléments moins fréquents. Quand les éléments sont équiprobables la longueur des codes doit être la même pour tous les éléments car il est plus avantageux de manipuler des codes à longueur fixe que des codes à longueur variable, étant donnée que l'efficacité des deux codes est presque la même.

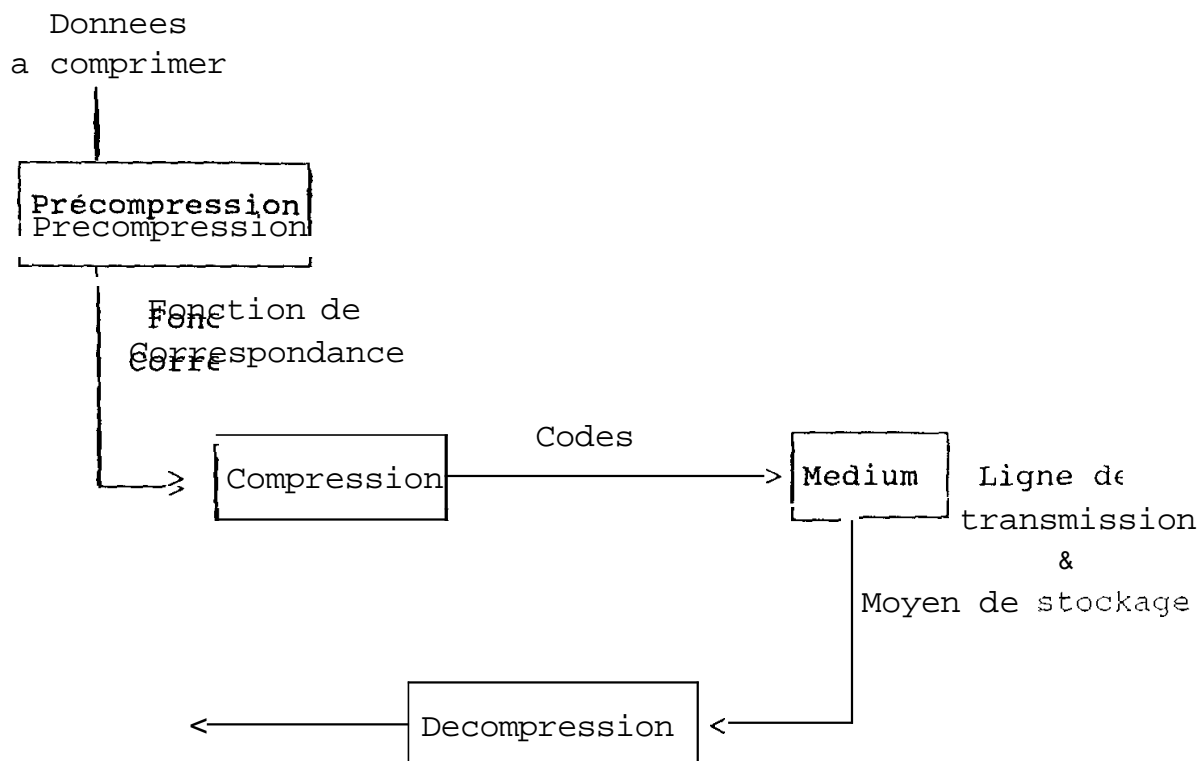
**II.2.3- Schéma de compression de données**

Ce schéma réalise les fonctions de codage, de decodage et de construction du code. Il doit être transparent à l'utilisateur, et est composé de trois principaux modules:

1- Le module de precompression, analyse le texte a compresser ( ou un echantillon representatif ) et a partir de ses caracteristiques, il calcule la fonction de correspondance, c'est a dire definit le code. Ce module peut être intégré au compresseur ou Btre independant de ce dernier.

2- Le module de compression (compresseur, encodeur) reconnait les messages delivres par la source et leur substitue dans la chaîne de sortie les codes equivalents.

3- Le module de decompression (decompresseur, decodeur) realise l'opération inverse en restituant les donnees initiales a partir des codes. Ceci peut Btre represente de la façon suivante:



#### 11.2.4- Parametres d'évaluation des TCD

Différents parametres d'évaluation sont utilisés pour juger de l'efficacite des techniques de compression de donnees et pouvoir les comparer. Les deux parametres les plus significatifs sont :

a- Le taux de compression atteint : la formule la plus utilisée est celle donnant le gain realise par rapport a la taille originale du texte, exprimé en pourcentage.

$$T = ( S - o ) / S$$

ou: S est la taille du fichier original.  
o est la taille du fichier comprime.

Le taux de compression differe d'une TCD a l'autre et varie generalement entre 40 et 90% [ALS 75].

b- Le temps de codage et de decodage

c- Les autres parametres sont: l'encombrement mémoire des programmes de compression/décompression et de la fonction de correspondance, la complexite de calcul des algorithmes ( des trois modules ), l'efficacité du code, etc.

#### II.3- Les redondances dans la representation des donnees

Toutes les TCD tentent d'améliorer la densité de stockage des donnees par l'élimination des redondances qui encombrant

la representation des donnees. Chaque TCD est conçue pour exploiter un certain type de redondance. Les TCD assez flexibles pour manipuler plusieurs types de redondances en même temps sont rares, et leurs resultats ne sont pas stables pour toutes les donnees [REG 81].

Les redondances les plus couramment rencontrees dans les fichiers et bases de donnees (donnees textuelles, numériques, alphanumeriques) sont:

- Une correlation significative existe entre des valeurs successives (dictionnaire de mots, thesaurus, ...).
- Longues sequences de caracteres repetees (blancs, zeros, ...).
- Certains caracteres se presentent avec une frequence bien supérieure a celle d'autres: non uniformite de la distribution statistique **des** caracteres.

#### 11.4- Classification des techniques de compression de donnees

Les TCD peuvent être classees de différentes manieres suivant le critère de classification choisi.

Dans ce paragraphe, nous donnerons les deux classifications les plus usuelles. La premiere sera basee sur l'élément de base de la codification et la deuxième sera basee sur la fonction de correspondance (dependance de la fonction de correspondance des donnees, sa construction et son fonctionnement).

### **II.4.1- Classification par élément de base de la codification**

Différents éléments de base peuvent être pris comme unité de codification, le choix de ceux-ci conditionne la performance et la complexité de la TCD:

#### **11.4.1.1- Les caractères**

Les TCD basées sur les caractères sont nombreuses et ont l'avantage d'être simples et non encombrantes du point de vue espace (table de Codage/Décodage très limitée). Cependant, l'utilisation du caractère comme unité de codage implique l'emploi de codes à longueur variable, ce qui est très contraignant.

Pour des données textuelles les performances de telles TCD sont limitées à moins de 50% [PIK 81].

#### **11.4.1.2- Les bigrammes**

Le codage se fait au niveau des caractères, mais l'alphabet est augmentée d'un certain nombre de bigrammes (paire de caractères les plus fréquents). Ces bigrammes, seront codés par les combinaisons (ASCII, EBCDIC) non utilisées par l'alphabet initial.

Ce type de TCD est simple à implémenter et réalise un taux de compression de 43 à 48% [WIS 87]. Il est important de signaler que dans cette méthode le rapport taux de compression/coût est intéressant. Par conséquent ces TCD sont préférées aux précédentes dans le cas des données textuelles.

Cependant, le problème réside dans le choix des bigrammes à coder, ce qui nécessite toute une analyse des textes à compresser.

Notons que ce principe peut être étendu à des trigrammes et des quadrigrammes, mais le gain obtenu est trop faible par rapport au coût, au fur et à mesure que la taille des chaînes de caractères augmente [WIS 87].

#### 11.4.1.3- Les fragments

Les TCD basées sur les fragments représentent une extension de celles basées sur les bigrammes. Dans celles-ci il y a création d'un dictionnaire de fragments (chaînes quelconques de caractères) les plus fréquents, et le texte sera codé comme une suite de fragments.

Ces techniques de compression donnent de bons résultats avec les données textuelles, mais leur coût peut être important. En effet, la taille du dictionnaire doit être assez grande pour permettre des taux de compression élevés, les algorithmes de précompression sont très complexes (après analyse du texte il faut définir l'ensemble optimal de fragments). Le même problème se rencontre avec les algorithmes de compression car ils doivent déterminer le découpage optimal du texte en fragments [MEC 89].

#### 11.4.1.4- Les mots

Pour les techniques basées sur les mots, le texte est codé comme une suite de codes de mots. L'inconvénient de ces TCD est la taille du dictionnaire de mots qui est très grande

(sauf si celui-ci ne contient que les mots ayant une fréquence assez élevée, auquel cas une deuxième TCD sera utilisée pour le codage de mots les moins fréquents [MEC 89] [HEA 72]).

#### 11.4.2- Classification suivant la fonction de correspondance

Suivant ce critère de classification, les TCD se répartissent en deux groupes [REG 81] [MEC 89], selon qu'elles dépendent ou non des données à comprimer.

##### 11.4.2.1- Les TCD sémantiquement dépendantes

Une TCD est sémantiquement dépendante si le codage qu'elle utilise dépend du contexte et de la sémantique des données à comprimer. Parmi ces TCD, nous avons la différenciation, le codage des répétitions [MEC 89], la compression des index [MEC 89] [WIS 86], etc.

##### ■ La différenciation

Une valeur est remplacée par sa différence avec une information de référence. La représentation des différences est moins encombrante que celle des valeurs originales, surtout quand elles sont très proches.

Avant compression	Après compression
Compression	Compression
Compilation	différenciation
Codage	différenciation

### - Codage des repetitions

Une longue sequence du même caractere est remplacee par un triplet contenant un indicateur de repetition, un compteur et le caractere repete.

#### Exemple :

Avant compression	Apres compression
122222	1*52
76999999	76*69

#### 11.4.2.2- Les TCD semantiquement independantes

Une TCD est dite semantiquement independante ou generale, si elle n'utilise aucune information liee aux donnees.

Cette generalite est limitee a une classe de donnees, ce qui assouplit la TCD pour plus d'efficacite dans l'implimentation de celle-ci.

Les TCD semantiquement independantes peuvent être utilisées avec n'importe quelles donnees, mais avec un degre d'efficacite variable. Parmi elles nous avons la TCD basee sur le code de HUFFMAN [MEC 89] [PEC 85], et celle basee sur le code de QOC (a Quasi-Optimum Code) [MEC 89] [FUM 86] [FUM 82].

### 11.5- Choix d'une TCD pour la compression des données bibliographiques

Les données contenues dans une base bibliographique se répartissent en trois types:

- **Données numériques:** Elles se trouvent dans les champs nombre de pages, numéro de volume, issn, isbn, etc.
- **Données textuelles:** Elles sont contenues dans les champs titre, résumé et mot-clé.
- **Données nontypees:** Cette classe regroupe toutes les rubriques qui ne contiennent ni du texte ni du numérique pur, tels que les champs auteurs, titre de conférence, lieu de conférence, etc.

Etant donnée la diversité des types de données constituant les éléments des fichiers bibliographiques, une TCD générale considérant le fichier comme des données homogènes risque de ne pas être efficace [MEC 89]. La meilleure alternative serait donc un schéma de compression dans lequel nous combinerons trois codages distincts, chacun adapté à un type de données (numérique, textuel et nontype) et pour lequel il pourrait donner les meilleurs résultats. Ce schéma doit respecter le profil suivant [MEC 89]:

1- La fonction de codage (3 fonctions) est calculée à partir d'un échantillon représentatif de la base documentaire.

2- Il doit réaliser le meilleur compromis entre un taux de compression très élevé et un temps de décompression très bas

(pour ne pas augmenter sensiblement le temps de reponse du systeme). La rapidité de la fonction de codage est secondaire car cette operation n'est effectuee que rarement (lors de la mise a jour mensuelle de la base).

3- L'espace mémoire nécessaire au stockage de la fonction de codage/décodage doit **être** négligeable devant le gain **realise**.

#### 11.5.1- Compression des donnees numeriques

Le codage des donnees numeriques est base sur les caractères [MEC 89]. L'expérience réalisée sur la base INSPEC donne un taux de compression egal a 65%.

#### 11.5.2- Compression des donnees nontypees

Cette classe de donnees est subdivisée elle-même en sept groupes:

- 1- Les dates
- 2- Les numeros de pages
- 3- Les numeros de volume
- 4- **Les champs d'identification**
- 5- Les champs de contrble
- 6- Les noms de personne
- 7- Les lieux et organisation

Par consequent, la technique préconisée est composee elle aussi de differents codages [MEC 89]. En effet, chaque groupe de rubriques prenant leurs valeurs dans le même domaine sera code par une technique adaptee au domaine. Les resultats obtenus pour ce type de donnees sont resumes dans le tableau **sui vant**:

Groupes de donnees	Taux de compression
<ul style="list-style-type: none"> <li>- Non typees</li> <li>  . Dates</li> <li>  . N° de pages</li> <li>  . N° de volumes</li> <li>  . Champs d'identification</li> <li>  . Champs de contrôle</li> <li>  . Noms de personnes</li> <li>  . Lieux, organisations</li> </ul>	<ul style="list-style-type: none"> <li>66,47%</li> <li>57,30%</li> <li>71,87%</li> <li>42,88%</li> <li>56,66%</li> <li>47,93%</li> <li>45,13%</li> </ul>

### 11.5.3- Compression de donnees textuelles

Les donnees textuelles sont contenues dans les champs: titre, resume et mots cles, et occupent la plus grande proportion du fichier bibliographique. Par consequent, si on veut avoir un taux de compression élevé pour le fichier bibliographique, il faut choisir un bon codage pour les donnees textuelles, c'est-a-dire un codage qui permet d'avoir un taux de compression eleve. Pour ce faire, deux techniques ont ete experimentees:

1- La premiere est basee sur les fragments, différentes variantes ont ete testees en jouant sur la taille du code, les valeurs 8, 9, 10, 11, 12 bits ont ete retenues ce qui a fixe le nombre de fragments codes a 256, 512, 1024, 2048, 4096.

2- La seconde est basee sur les mots. Dans cette technique les mots a coder sont tries par ordre décroissant des frequences et chaque mot recevra comme code son rang dans l'ordre établi.

L'implémentation de ces TCD a donne les resultats resumes dans le tableau ci-dessous:

TCD	Taux de comp. (%)	Temps de codage ms/car	Temps de decodage ms/car	Encombrement mem Ko	
				codage	décodage
F256	45,93	0,053	0,025	3	2
F512	48,54	0,122	0,105		
F1024	53,34	0,120	0,092		
F2048	53,60	0,124	0,091	26	26
F4096	56,28	0,128	0,081	52,5	50
Mot	71,23	0,392	0,146	93	96

**Remarques:**

1- L'encombrement memoire est d'autant plus important que le taux de compression est grand. Aussi nous remarquons que le codage des fragments sur un octet (**F256**) est le moins encombrant du point de vue espace memoire.

2- Pour le temps de decompression la TCD **F256** est la plus efficace car elle ne necessite pas la manipulation de bits. Le temps de decompression de la TCD basee sur les mots est le plus important car elle travaille sur des codes a longueur variable, les autres TCD ont des temps de decodage tres proches.

3- Le temps de compression est tres important dans la TCD basee sur les mots car le codage implique une recherche dans une table. Ce temps est le même pour les autres TCD sauf la F256 qui reste la plus rapide.

4- Le taux de compression realise est nettement plus important pour la TCD basee sur les mots. En effet, ce type de TCD offre un taux de compression appreciable (premier critere de choix d'une TCD) par rapport a celles basees sur les fragments, et un taux de decodage acceptable (il peut être optimisé par la programmation du decodeur en assembleur). C'est pour toutes ces raisons que nous l'avons choisie, dans le cadre de notre systeme, pour la compression des donnees textuelles [MEC 89].

#### 11.6- Integration de la technique de compression de donnees au systeme SIBADOC

Le Systeme d'Interrogation des Bases DOCUMENTAIRES SIBADOC a ete conçu et developpe au CERIST dans le but de prendre en charge la gestion et l'accès a l'information documentaire contenue dans des bases de donnees bibliographiques internationales.

SIBADOC a ete developpe sur vax-11/785 et est constitué de deux principaux modules [BEN 89]:

- le module de creation et de mise a jour des bases documentaires,
- le module d'interrogation pour la recherche

Une presentation detaillee de SIBADOC est donnée en annexe A.

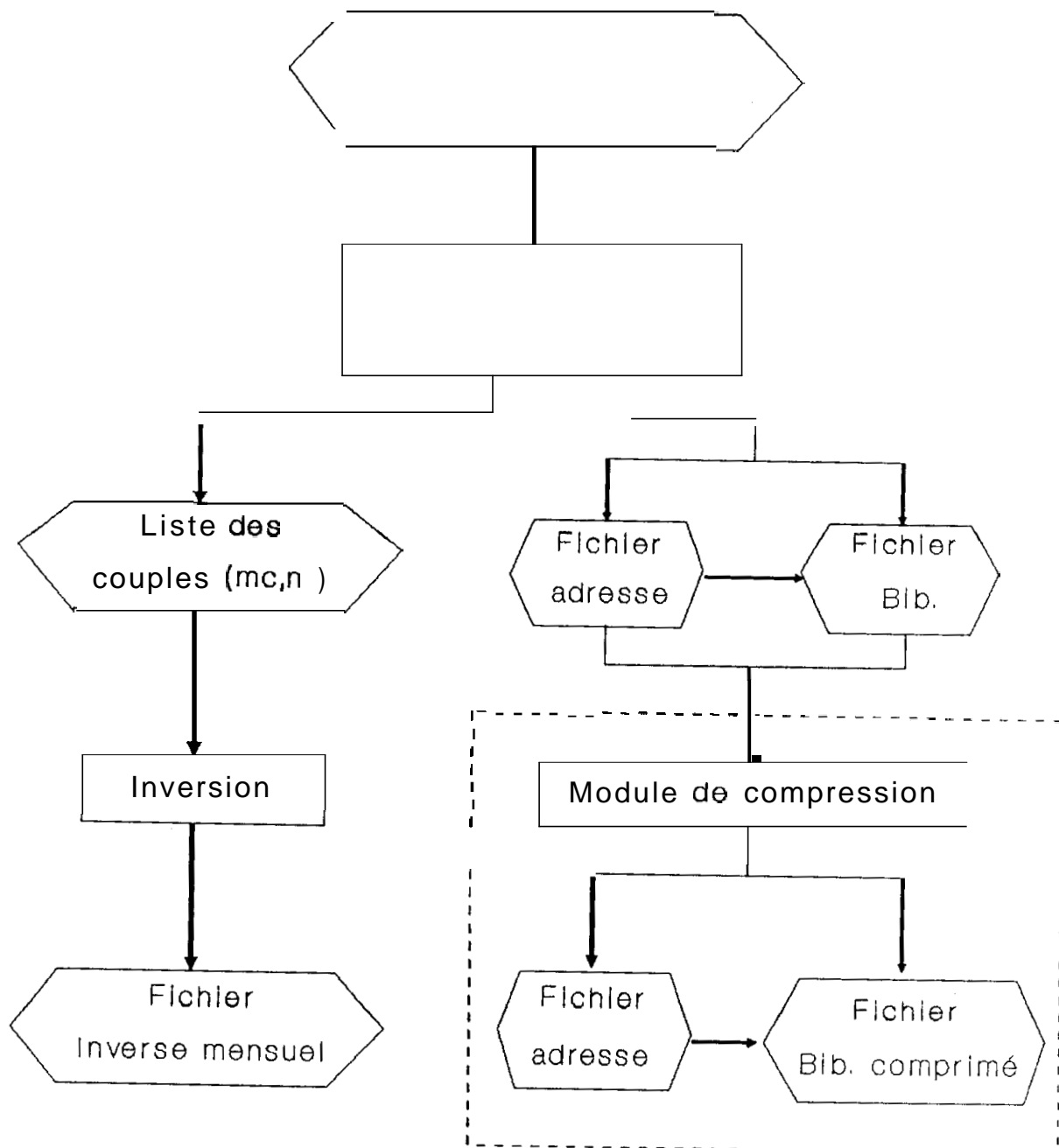
Après avoir choisi la TCD appropriée aux données bibliographiques, cette dernière a été intégrée à titre expérimental au système SIBADOC dans le cadre d'un projet de fin d'études [MEC 89]. Le gain réalisé par le module de compression est satisfaisant, puisque ce dernier a permis un taux de compression de 64% sur la base INSPEC.

Nous présentons dans ce qui suit la localisation dans le système SIBADOC des actions de compression et de décompression. Cette localisation est très simplifiée étant donnée la structure modulaire du système SIBADOC.

#### **II.6.1- La compression dans SIBADOC**

Le module de compression est intégré au module de création et de mise à jour et il n'intervient qu'à la fin de cette opération. Cet algorithme réduit considérablement la taille des notices et met à jour leurs adresses dans le fichier adresse correspondant.

Le schéma d'exécution de ce module devient :

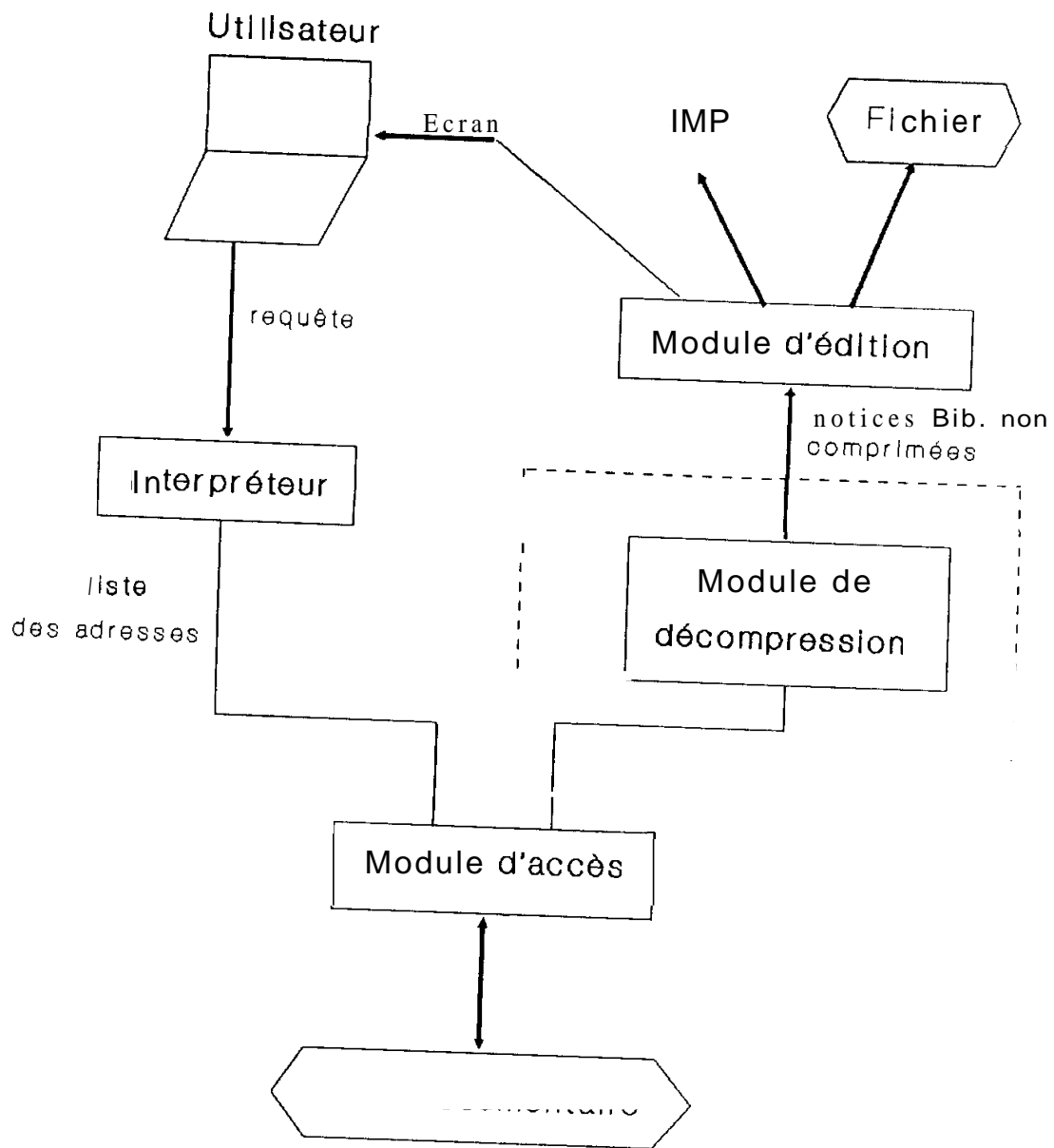


**Module de creation et mise a jour  
de la base avec compression de donnees.**

#### 11.6.2- La décompression dans SIBADOC

Le module de décompression est intégré au module d'interrogation de la base, il sera exécuté immédiatement avant l'édition des notices demandées par l'utilisateur.

Le processus d'interrogation se schématisera alors comme suit:



**Module d'interrogation de la base comprimée**

Nous terminons la presentation de l'intégration de la technique de compression de donnees au systeme SIBADOC en donnant des resultats par rubriques obtenus dans la compression du fichier du mois de Janvier de l'année 1989, de la base INSPEC en exploitation sous SIBADOC, ainsi que les resultats de la compression de toute la base INSPEC (6 mois sont en exploitation).

Groupes de donnees	Taux de compression
- Numériques	65,57%
- Textuelles	71,21%
- Non typées	
. Dates	66,47%
. N° de pages	57,30%
. N° de volumes	71,87%
. Champs d'identification	42,88%
. Champs de contrble	56,66%
. Noms de personnes	47,93%
. Lieux, organisations	45,13%
Moyenne de compression	64,10%

Resultats par rubrique de la compression  
de INSPEC0189

Fichier INSPEC	Taille du fichier avant compression	Taille du fichier apres compression	Taux de compress.
Janvier	31066	10732	65,45%
Février	29512	10254	65,26%
Mars	20432	7136	65,07%
Avril	35640	1251%	64,88%
Mai	34334	11982	65,10%
Juin	42312	14734	65,18%

Resultats de compression de la base INSPEC

Après cette étude des méthodes de compression de données existantes, nous allons passer à la détermination d'une méthode d'accès optimale et qui s'adapte aux données comprimées. Ceci fera l'objet du chapitre suivant.

*CHAPITRE III*  
**TECHNIQUES D'ACCES**

Dans la littérature, diverses techniques d'accès aux données textuelles ont été proposées. L'objectif de ce chapitre est d'effectuer une étude exhaustive des différentes techniques d'accès et de proposer une technique optimale d'accès aux données volumineuses.

### 111.1 Description des différentes techniques d'accès

Les techniques d'accès aux données se focalisent toutes autour des grandes classes suivantes :

#### III.1.1 - Recherche Séquentielle (Full Text Scanning)

Etant donné un échantillon de recherche, la base de données sera entièrement consultée jusqu'à ce que les documents en question soient retrouvés. Cette méthode ne nécessite pas d'espace overhead et ne demande pas beaucoup d'effort pour l'insertion, ni pour la mise à jour [FAL 84]. Cependant, le principal inconvénient est le temps de recherche. Bien qu'il existe quelques algorithmes de recherche d'une chaîne de caractères qui sont rapides [BOY 77] [KNU 77] [AHO 75] [MAL 71], la recherche dans une grande masse de données peut nécessiter un temps considérable [HOL 79].

#### III.1.2- Inversion

Plusieurs systèmes, tels que STAIRS [RAB 4] [IAE 82], LEXIS et SIBADOC ont adopté cette approche.

Cette methode utilise un index. Une entree de cet index consiste en un mot (ou theme ou concept), suivi d'une liste de pointeurs. Ces pointeurs donnent l'adresse des documents qui contiennent ce mot.

Le principal avantage de la methode est sa rapidité de recherche. Cependant, elle peut necessiter un large espace overhead de stockage pour l'index. Cet espace peut aller jusqu'à 3 fois la taille du fichier contenant les donnees si les informations de localisation des mots sont egalement pris en compte [FAL 84] [FAL 85].

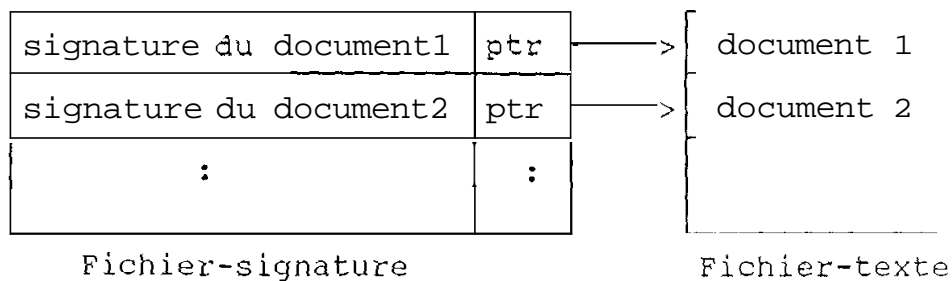
### 111.1.3 - Regroupement (Clustering)

Dans cette technique d'accès, les documents similaires sont regroupes ensemble pour former des groupes, appeles clusters, conformement au principe de regroupement qui stipule que les documents tendant a répondre aux mêmes questions sont étroitement lies, c'est-à-dire qu'ils appartiennent au même groupe [EAS 89] [BEL 87].

Il semble difficile de comparer cette methode avec les precedentes, car le regroupement se fait selon le sens des requêtes. Par exemple, "donnez-moi les documents qui sont en rapport avec la recherche d'information (même si les documents ne contiennent pas les deux mots 'Recherche' et 'Information'). Il semble que l'espace overhead est petit, de l'ordre de  $\log n$ , ou  $n$  est le nombre total de documents. De même, le temps de recherche est de l'ordre de  $\log n$ . Par contre, cette methode semble ne pas pouvoir gérer facilement les insertions [FAL 84].

## 111.1.4 - Fichier-signature (signature File)

De par le monde, de nombreux travaux concernant les signatures et l'accès utilisant les fichier-signatures ont été réalisés [RAB 84] [ROB 79] [SAC 87] [KUO 82] [TAV 92]. Dans cette technique d'accès, les documents sont stockés séquentiellement dans un fichier texte. A Chaque document est associé un mot-code appelé signature du document, obtenu en appliquant une fonction de hachage à ce dernier. Les mots-code sont de même taille (F bits) et sont stockés séquentiellement dans un fichier **appelle** fichier-signature.



Quand une requête arrive, au lieu de consulter chaque document, on consultera la forme condensée de ces documents, c'est-à-dire le fichier-signature. De ce fait, un grand nombre de documents qui ne satisfont pas à la **requête seront écartés, et par conséquent, le texte d'un document ne sera vraiment consulté que s'il y a une forte probabilité d'y retrouver l'information désirée.** Le fichier-signature peut être alors considéré comme un résumé ou des mots clés.

Cette methode est plus rapide que celle du 'Full text scanning', mais elle est supposee être plus lente que l'Inversion [FAL 84]. Cependant, la methode des signatures necessite moins d'espace overhead que celle de l'Inversion et peut gerer facilement les insertions [CHR 84] [JOH 69] [TSI 82].

Cette breve discussion sur les methodes de recherche textuelle laisse voir qu'aucune de ces methodes n'est clairement meilleure que les autres. Par consequent, nous allons proposer deux autres techniques d'accès qui s'intègrent tres bien aux donnees comprimees et dont l'objectif est d'accélérer la recherche d'information :

#### 111.1.5 - Methode d'accès utilisant un Fichier-signature

Dans cette technique, des listes inversees sur des partitions comprimees de la base de donnees seront utilisées.

L'inversion se fait au niveau des symboles (mots ou fragments de mots), utilises pour la compression de donnees [GOY 85]. Cette methode s'intéresse donc aux symboles qui ont la particularité d'être tres frequents, etant donne qu'il existe une forte correlation entre les mots sélectionnés par les indexeurs et ceux qui apparaissent de maniere frequente dans un texte [KEV 83]. La base de donnees consistera donc en des pages comprimees auxquelles sont associés des index (un index par page).

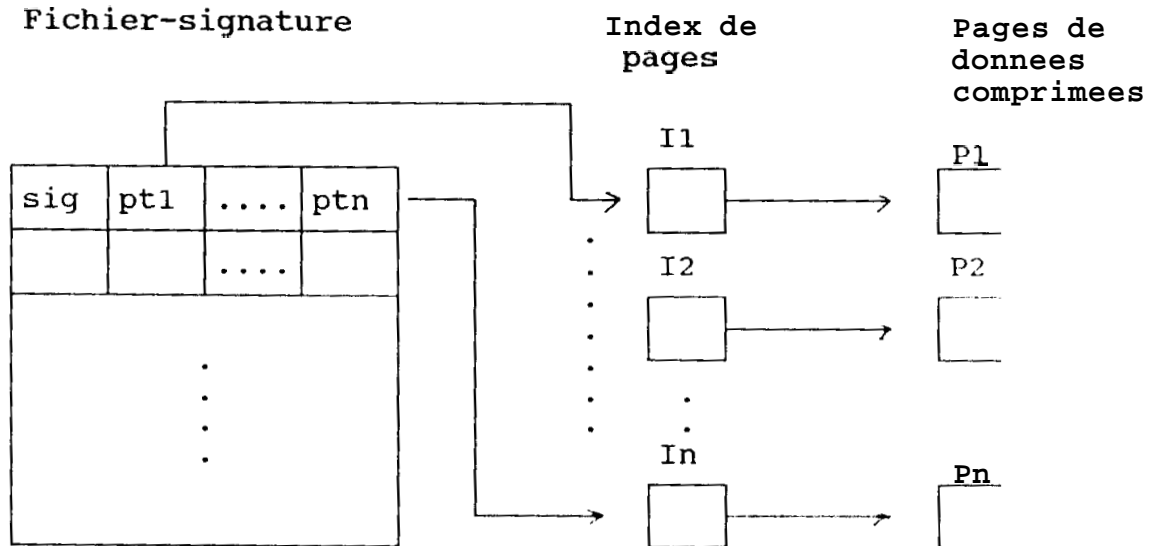
Pour rechercher une chaîne 'S', celle-ci est eventuellement divisée en plusieurs symboles 'S1''S2'...'Sm'

tel que chaque 'Si' appartienne a l'ensemble des symboles utilises pour la compression. Ces 'Si' sont utilises pour determiner les index en question.

Pour augmenter les performances de la recherche, en plus des index de pages, un fichier-signature contenant non pas des symboles, mais les mots-code (ou signatures) de ces symboles, sera utilise [GOY 85] (voir fig 111.1). Ces signatures sont de même taille (f bits) et sont obtenues en appliquant une fonction de hachage a l'ensemble des symboles contenus dans les differents index de pages.

Le fichier-signature contiendra donc les signatures de tous les symboles contenus dans les differents index de pages (une signature par symbole), ainsi que des pointeurs vers les differents index qui contiennent ces symboles.

Ainsi, au lieu de parcourir tous les index de pages pour determiner ceux qui contiennent un symbole donne, on consultera le fichier-signature, et a l'aide des pointeurs, on selectionnera un sous-ensemble d'index de pages pour la recherche, ce qui augmente les performances de la technique d'accès.



signature d'un symbole	pointeur vers index1	...	pointeur vers indexn
------------------------	----------------------	-----	----------------------

#### Une entree d'un indexi

symbole j	adresse de la notice contenant le symbole j dans la page i
-----------	------------------------------------------------------------

Fig 111.1: Acces en utilisant un fichier-signature

### 111.1.6 - Méthode d'accès utilisant un Q-arbre

#### III.1.6.1 - Définition d'un Q-arbre

Un Q-arbre est un arbre équilibré comme les B-arbres [BAY 72]. En effet, tous les chemins de la racine aux feuilles sont de même longueur. La différence entre un Q-arbre et un B-arbre est que la profondeur d'un Q-arbre est fixée à 1, et la taille de ses pages est variable (la taille de la racine est différente de celle des feuilles) [FUM 86].

La technique du Q-arbre est applicable pour n'importe quel index qui satisfait les conditions suivantes:

a) A la construction de l'index, il existe déjà un grand ensemble de mots. Ces mots sont appelés mots primaires.

b) Les mots qui vont être ajoutés à l'index seront appelés mots secondaires.

c) Aucun mot primaire ne sera supprimé, car la suppression d'un mot nécessite la reorganisation des feuilles du Q-arbre. De plus, les Q-arbres prennent en compte le fait que les mots clés sont rarement supprimés dans la plupart des systèmes de recherche d'information.

d) Le nombre de tous les mots (primaires et secondaires) est supérieur à plusieurs millions.

### 111.1.6.2 - Description du Q-arbre

Un element de l'index est un couple  $(w, q)$ .

ou  $w$  : est le mot qui constitue la cle de recherche.

$q$  : est un pointeur vers les informations qui sont associees a  $w$ .

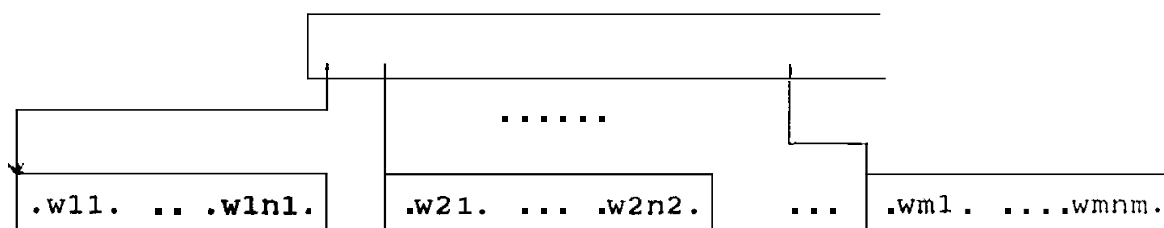
Toutes les feuilles du Q-arbre ont la même taille et sont stockees en memoire secondaire.

La racine reside en memoire centrale et peut ne pas avoir la même taille qu'une feuille. En fait, la taille de la racine depend de l'espace memoire dont le systeme peut disposer.

Les mots dans la racine et les feuilles sont tries dans l'ordre alphabetique.

Tous les mots secondaires sont stockes dans les feuilles du Q-arbre.

Contrairement aux feuilles, la racine n'a pas d'espace pour les ajouts. De plus, la racine ne contient que des mots primaires et vu les conditions d'application de la technique du Q-arbre (cond. c), ces mots ne seront pas supprimes.



**Arbre équilibré de profondeur 1**

### 111.2 - Critères de choix de la technique d'accès

C'est l'organisation du fichier inverse de mots clés qui détermine l'efficacité d'un système de recherche d'information.

Concernant la technique d'accès utilisant un fichier-signature à la place de ce dernier, nous aurions pu utiliser la technique d'indexation traditionnelle; c'est-à-dire utiliser un autre niveau d'indexation et qui sera un index global contenant tous les symboles au lieu de leurs signatures, mais l'apport de la méthode des signatures en gain d'espace mémoire est très important en comparaison avec celui des techniques d'indexation usuelles. Néanmoins, cette technique présente toujours des inconvénients du fait qu'elle nécessite un grand espace overhead dû à l'utilisation de deux niveaux d'indexation: le fichier-signature et les index de pages. De plus, l'utilisation d'une fonction de hachage pour l'obtention des signatures peut engendrer des problèmes de collision, ce qui augmente le nombre de réponses non pertinentes par rapport aux questions et par conséquent, influe sur l'efficacité du système.

D'autre part, le facteur principal qui détermine l'efficacité d'un arbre équilibré est sa profondeur et elle est minimale dans un Q-arbre. La structure du Q-arbre semble donc être l'organisation la plus appropriée qui permet d'avoir un index rapide [MAT81].

Comme l'index de mots clés satisfait les conditions d'application de la technique du Q-arbre, alors, c'est cette structure qui a été adoptée dans notre système pour l'index de mots clés.

La solution retenue en utilisant un Q-arbre, est donc de placer les mots les plus fréquents dans la racine du Q-arbre, ce qui diminue le nombre d'accès disque pour la recherche, et de placer le reste des mots dans les feuilles du Q-arbre.

L'ensemble des mots de départ, qu'ils soient stockés dans la racine (mots les plus fréquents) ou dans les feuilles, seront appelés alors mots primaires.

La détermination de l'ensemble des mots utilisés par le Q-arbre se fait par une méthode d'indexation dont l'étude et le choix feront l'objet du chapitre suivant.

## *CHAPITRE IV*

*INDEXATION OU SELECTION DE MOTS CLES*

L'indexation consiste à trouver ou à sélectionner le mot vedette, appelé aussi mot **clé**, qui est censé traduire le mieux le sujet principal dont traite chaque document [LON 80].

L'indexation est un art qui réclame beaucoup de qualités, souvent contradictoires de la part de l'indexeur humain [CHA 90]: compréhension du domaine, compréhension de la langue, esprit analytico-synthétique, etc.

Les produits de l'indexation sont des index, c'est-à-dire des listes de termes significatifs, qui servent à retrouver les informations et à les sélectionner afin de répondre aux besoins des utilisateurs.

Comme toutes les activités d'information, l'indexation est d'abord un instrument de travail. Elle doit permettre de retrouver les informations utiles pour des catégories bien précises d'utilisateurs effectuant des travaux déterminés. L'indexation doit donc être définie de telle sorte que ces besoins puissent être satisfaits au moindre coût.

De plus, la qualification et les effectifs du personnel disponible, le volume d'informations à traiter et leur nature, le système de stockage et de recherche de l'information, la nature de la matérialité des produits et des services de diffusion de l'information, les moyens financiers disponibles imposent également des limites dans le choix d'une méthode d'indexation.

#### **IV.1- Les méthodes d'indexation**

Du point de vue de sa finalité, l'indexation est destinée à permettre une recherche retrospective efficace des informations contenues dans un fonds documentaire.

L'indexation est donc la tâche la plus importante dans un système de recherche documentaire. En effet, c'est elle qui conditionne l'efficacité de ce dernier. Une mauvaise indexation ou une indexation insuffisante représente 90% des causes essentielles de l'apparition, lors d'une recherche, de bruits ou de silences [CHA 80].

Le bruit est le taux de documents parasites qui ne répondent pas à la question posée et qui sont extraits lors des opérations de sélection. Si pour 100 documents obtenus en réponse, 20 sont jugés ne pas répondre à la question, l'on dit que le bruit est de 20%.

Le silence est le taux de documents pertinents existant dans la base documentaire qui devraient répondre à une question et qui n'ont pas été sélectionnés lors de l'interrogation. Le silence est dû à l'emploi de caractéristiques non concordantes dans la formulation de la question et dans l'indexation du document. Un système donnant un taux de silence d'environ 30% est généralement considéré comme satisfaisant.

#### **IV.1.1- Indexation sans thesaurus et indexation avec thesaurus**

##### **IV.1.1.1- Définition d'un langage documentaire**

Un langage documentaire est un ensemble de termes utilisés pour représenter le contenu des documents [CHA 73] [RIV 90].

Les éléments constitutifs d'un langage documentaire sont les suivants [GUI 90]:

- Des mots servant à décrire les informations, les descripteurs. Ils sont tirés du langage naturel, réduits à une

forme grammaticale unique et invariable (généralement le substantif singulier). Ils peuvent être simples ou composés.

- Des mots du langage naturel ayant des relations avec les descripteurs et qui sont repertoriés avec un renvoi au descripteur correspondant. Ces mots ne peuvent être utilisés pour décrire les informations. On les appelle des non-descripteurs.

- Des relations entre les descripteurs: relations hiérarchiques, d'équivalence ou de voisinage, etc. Elles permettent de regrouper les notions sous un seul terme, d'élargir ou, au contraire, de préciser une recherche.

#### **IV.1.1.2- Définition d'un thesaurus**

Le thesaurus représente un mode d'organisation d'un langage documentaire. C'est un répertoire de mots qui ont entre eux des relations (hiérarchique, voisinage ou équivalence, synonymie, ...), dans un domaine particulier de connaissance.

Le thesaurus est utilisé par les indexeurs au moment où les concepts d'indexation doivent être transformés en descripteurs [GUI 89].

Les concepts font partie du langage naturel, les descripteurs du langage contrôlé du thesaurus.

Selon qu'il dispose ou non d'un thesaurus, l'indexeur aura recours à des procédés différents:

#### IV.1.1.3- Indexation sans thesaurus

Cette methode est frequemment appelee methode d'indexation en "langage naturel" car le rôle du mot cle (ou descripteur) est devolu, non par comparaison a une liste de mots cles existants deja, mais a tout candidat mot cle.

Cette methode offre une grande souplesse et une simplicité de mise en oeuvre, mais presente les inconvenients suivants:

- Tres coûteuse (en espace mémoire a cause de la grande taille des fichiers).

- On peut aboutir rapidement a une dispersion catastrophique du lexique. Pour y remedier, il faut au moins:

- \* eliminer les mots non significatifs appelés aussi mots vides (articles, prepositions, ...).

- \* réduire a une seule réalisation morphologique tous les mots presentant des variantes paradigmatisques (nombre, genre, conjugaison, déclinaison).

#### IV.1.1.4- Indexation avec thesaurus

Dans cette methode, l'indexeur recherche si les mots cles qu'il a repéré existent dans le thesaurus, sinon il les traite comme des candidats descripteurs en exerçant les dcux contrôles suivants:

- 1- Examen de la liste des candidats descripteurs et refus de certains d'entre eux (soit par mise a jour de l'ensemble des

mots non significatifs, soit par rejet pur et simple, soit par insertion dans une liste d'attente), car l'apparition du candidat descripteur est estimée comme étant tout à fait accidentelle ou non significative dans le corpus donné.

2- Examen de la liste des candidats descripteurs et remplacement de certains d'entre eux par des descripteurs estimés avoir une signification identique ou une signification proche et ceci en établissant les relations classiques (hiérarchie, synonymie, voisinage, association, ...) entre descripteurs.

L'indexation avec thesaurus nous conduit donc à une liste structurée de descripteurs. Plus cette liste est structurée moins il est exigé d'aptitudes de l'utilisateur, qui peut être non spécialisé et plus le système de recherche documentaire sera performant en qualité [DEW 81] [DEW 85].

#### **IV.1.2- Indexation manuelle et indexation automatique**

Il existe deux principales méthodes d'indexation liées d'une part à l'importance du fonds documentaire et d'autre part à la disponibilité d'un équipement informatique:

- 1- Indexation manuelle
- 2- Indexation automatique.

##### **IV.1.2.1- Indexation manuelle**

L'indexation manuelle est une indexation faite par des indexeurs humains, spécialistes des domaines abordés dans les documents.

L'indexation manuelle est une solution satisfaisante dans un contexte particulier [TIG 92]:

- 1- Nombre de documents peu élevé.
- 2- Textes ne se trouvant pas sur support informatique.

Cependant, cette methode d'indexation s'avère inadéquate dans les grands centres de documentation où le volume de l'information est tres important. De plus, ce type d'indexation presente les inconvenients suivants [CHA 89] [TIG 92] :

- 1- Rigidité de l'analyse et des concepts qui y découlent.  
En effet:
  - \* Indexation figee
  - \* Mise a jour impossible, surtout dans des domaines qui evoluent tres rapidement tels que l'intelligence artificielle ==> effet de décalage entre le thesaurus et le contenu des documents.
  - \* Manque de souplesse face a l'évolution dynamique des besoins des utilisateurs.

2- Qualite de l'indexation réduite a cause:

- \* De l'effet de variabilité ou de non uniformité : deux indexeurs choisissent rarement exactement les mêmes termes pour indexer un même document. On observe, parfois, l'absence de certains descripteurs que l'indexeur aurait dû retenir pour rendre compte du contenu.

\* Du bruit (informations inutiles) et du silence (perte d'informations) documentaire qui sont conséquence de l'absence d'uniformité dans une indexation manuelle.

Les difficultés liées au processus d'indexation manuelle ainsi que les volumes très importants d'informations à traiter obligent et justifient l'apparition de systèmes d'indexation automatique de documents.

#### IV.1.2.2- Indexation automatique

L'indexation automatique [SAL 89] consiste à faire reconnaître par l'ordinateur des mots figurant dans le titre, le **résumé** ou le texte lui-même d'un document.

Les termes reconnus sont incorporés dans le fichier de recherche et servent à retrouver le document.

L'indexation automatique suppose que l'on introduise dans l'ordinateur sinon tout le texte, du moins le titre et un résumé d'auteur. L'ordinateur utilise alors différentes approches pour identifier les termes significatifs.

##### IV.1.2.2.1- Approche statistique

Dans cette approche on procède à des analyses statistiques sur un échantillon de textes pour déterminer la fréquence d'apparition des mots. Le petit nombre qui apparaît constamment **est** considéré comme sans signification. Un second groupe comprend des mots apparaissant moins souvent et considérés comme significatifs.

Un troisième groupe comprend des mots qui apparaissent très rarement et sont considérés comme trop spécifiques. A partir de là, on peut indiquer à l'ordinateur quelle doit être la fréquence des mots à retenir pour l'indexation.

La faiblesse de cette méthode, qui n'est applicable que dans un domaine donné, tient à ce qu'elle considère chaque mot isolément et permet seulement de dire qu'un concept ou un objet figure dans le document, sans préciser son rôle. Cette méthode ne tient pas compte, par exemple, du problème **pose** par la synonymie. Une même information peut apparaître plusieurs fois dans un document, mais sous des formes différentes. Cela sera ignoré dans le calcul, fait par l'ordinateur, de fréquence d'apparition d'un mot. Des dictionnaires de synonymie peuvent être établis mais ils alourdissent le système [GUI 90]. Ainsi, le calcul de la fréquence s'effectue sur le concept et non sur le mot.

#### IV.1.2.2.2- Approche linguistique

Vue l'insuffisance de l'approche statistique et la nature des informations traitées (langue naturelle) dans les systèmes de recherche documentaire, un modèle linguistique pour déterminer les concepts importants d'un document, s'impose naturellement.

Les méthodes linguistiques introduisent des traitements morpho-syntaxiques pour arriver à une reconnaissance des structures les plus significatives. L'indexation s'effectue alors en trois grandes phases [LAL 87] [MET 88] [COY 67] [ROU 87]:

- La première consiste à déterminer les unités lexicales du document à analyser. C'est l'analyse morphologique.
- La deuxième consiste dans la reconnaissance sémantique des termes retenus. Cela revient à classer ces derniers en différents groupes, selon qu'ils soient significatifs ou non. C'est l'analyse sémantique.
- La troisième et dernière phase consiste à établir des liens entre les descripteurs pris deux à deux issus des deux premières étapes. C'est l'analyse syntaxique.

L'indexation automatique vise donc à déterminer les concepts importants d'un document en examinant la suite de caractères qui forment le texte. Pour ce faire, il faut:

1- Bien connaître les mécanismes du langage qui permettent de passer de la forme au sens.

2- Être capable de simuler ce processus sur ordinateur.

Cependant, ces deux conditions ne sont pas remplies car la question qui se pose est de savoir comment simuler un processus, qui permet de passer de la forme au sens, sur ordinateur, si on ignore les mécanismes fondamentaux de ce processus. Ceci a amené certains chercheurs à dire que la véritable indexation automatique n'est qu'à l'état de projet de recherche [MET 88].

#### **IV.2- Critères de choix de la méthode d'indexation**

Les problèmes essentiels de la documentation sont des problèmes sémantiques. L'analyse et l'indexation des documents

en vue de leur mémorisation forment un goulot d'étranglement dans tous les systèmes documentaires (la production documentaire augmente plus vite que le nombre d'indexeurs qualifiés) et cet obstacle doit être surmonté pour que le transfert de l'information s'effectue dans de bonnes conditions. Ainsi différentes voies faisant appel aux ressources de l'informatique ont été explorées pour tenter de suppléer à l'indexation humaine (manuelle) et la remplacer:

- d'une part avec **les systèmes dits en texte intégral ou en "full text"**, systèmes supprimant la phase d'indexation et reposant sur une stratégie de recherche complexe.
- d'autre part avec les méthodes d'indexation automatique en amont de systèmes documentaires devenus classiques.

#### IV.2.1- Les systèmes en texte intégral ou systèmes en "full text"

Ces systèmes dits en "full text" ou encore "free text" ne font appel à aucune indexation, le texte intégral du document étant mémorisé et la recherche portant sur l'ensemble du texte et non plus sur les descripteurs [CHA 73] [CHI 86] [SAL 86b].

Généralement les systèmes en texte intégral portent non pas sur le document **lui-même**, sauf s'il est court (dépêche de presse, coupures de presse, articles de loi, jurisprudence, etc), mais sur le résumé plus ou moins développé du document.

Dans de tels systèmes utilisant le langage naturel, la recherche doit porter sur l'ensemble du texte mémorisé, aucune sélection de termes ni aucune indexation n'ayant été opérées.

Un element important qui rend ces systemes tres interessants est le fait qu'ils ne comportent aucune indexation; ce qui permet d'éviter d'avoir recours a des indexeurs humains dont le travail n'est jamais entierement satisfaisant (& IV.1.2.1). Cependant, une evaluation détaillée de ce type de systemes, faite dans [BLA 85] et [MAR 90], montre que ces systemes ne sont performants que pour des bases de donnees non volumineuses. De plus, ces systemes ne reposant sur aucune indexation, utilisent des strategies de recherche tres complexes telles que [CHA 73]:

- Utilisation de masques sur les mots du texte.

Par exemple, en écrivant dans l'équation de recherche DOCUMENTS, la recherche portera sur tous les mots ayant cette racine. Il sera alors possible de selectionner les documents ayant l'un des termes:

- \* DOCUMENTS
- \* DOCUMENTATION
- \* DOCUMENTAIRE
- \* DOCUMENTALISTE

- Utilisation d'une logique de recherche tres developpee: possibilités d'utiliser de nombreux extenseurs ou restricteurs (plus petit que, plus grand que, compris entre, adjacent, avec, etc).

**Exemple:**

La formulation TRAITEMENT AVEC DOCUMENTATION permettra de selectionner les documents ayant une phrase telle que: Le TRAITEMENT par des methodes informatiques de la DOCUMENTATION scientifique.

- Possibilité d'interrogation sur des zones déterminées, et ceci en utilisant des indicateurs de zones (zone auteur, zone bibliographie, etc).

Exemple:

En écrivant une question:

CONTROLE ZONE **BIBLIOGRAPHIE** SALTON

On fera apparaître tous les documents citant SALTON dans leurs références bibliographiques.

#### **IV.2.2- Méthodes d'indexation automatique**

Comme nous l'avons vu au chapitre IV, l'indexation automatique se résume en deux principales méthodes [SAL 70]:

- méthode statistique
- méthode linguistique

##### **IV.2.2.1- Méthode statistique**

Les méthodes d'indexation automatique les plus anciennes sont les méthodes statistiques [ROU 87] [CHA 73]. Avec ces méthodes, les mots les plus représentatifs sont extraits par un calcul de fréquence.

L'un des inconvénients majeurs de ces méthodes réside dans la non prise en compte, pour le calcul de la fréquence, des synonymes. Plusieurs méthodes sont utilisées pour tenter de remédier à cet inconvénient. La plus simple consiste à dresser un dictionnaire de termes synonymes, de manière à ce que le calcul de la fréquence soit effectué sur le concept et non sur le mot.

Pour affiner ces methodes statistiques, envers lesquelles de nombreuses critiques ont ete formulees, d'autres elements sont parfois pris en compte, en affectant des poids aux termes [LAM 82] [COO 78] ou en utilisant des operateurs linguistiques [CHA 73]. Ainsi il a ete propose de retenir en priorité comme termes significatifs ceux qui sont precedes d'expressions comme: en resume, en conclusion, etc. Toutefois, le probleme crucial et qui restera toujours pose, est la determination du seuil **de frequence** a partir duquel nous deciderons si un mot est significatif ou non [ROU 87].

#### **IV.2.2.2- Methode linguistique**

Devant les insuffisances des methodes statistiques, d'autres methodes ont ete experimentees, faisant appel aux ressources de la linguistique. Ces methodes, necessitent au depart un long travail de preparation d'outils d'analyses, constitution de dictionnaires, construction de reseaux scmantiques, elaboration de grammaires, etc. Or ceci fera l'objet de tout un projet de recherche, ce qui nous ai impossible de traiter en detail dans le cadre de cette these.

#### **IV.3- Conclusion**

Apres cette longue discussion sur les methodes d'indexation qui laisse voir que la meilleure approche pour indexer des documents est l'approche linguistique, mais qui n'est pas realisable dans le cadre de cette these, la methode qui a ete choisie pour indexer notre corpus documentaire, se resume dans un premier temps a extraire les termes du fonds documentaire et a eliminer ceux qui sont non significatifs par consultation d'un dictionnaire de mots vides.

Les mots vides representent l'ensemble des mots qui ne sont pas significatifs pour une recherche retrospective et ne figurent donc pas dans le fichier inverse d'une referotheque [DEW 81] [DEW 89] [GET 87] [MER 84].

Le choix des mots vides est d'une importance directe avec la taille de l'index ainsi qu'avec la qualité du systeme.

Examinons quelques ensembles de mots vides :

1- Ensemble des 'n' plus frequents mots.

'n' etant un parametre du systeme qui doit être choisi de sorte que la taille du fichier inverse soit réduite et que l'efficacité de recherche du systeme ne soit pas deterioree.

2- Ensemble des articles, prepositions, conjonctions, pronoms et adverbés.

Théoriquement, en utilisant cet ensemble de mots vides, la taille du fichier inverse peut être réduite de 60% (approximativement) avec 2% de perte en efficacité de temps de recherche [FUM 86].

Pour une telle methode de selection de mots vides, une reduction de 40% est au minimum garantie [FUM 86].

# *CHAPITRE V*

*SYSTEMES D'ACCES AUX DONNEES VOLUMINEUSES*

Dans ce chapitre nous presentons le systeme d'accès aux donnees volumineuses que nous avons conçu, son integration au systeme SIBADOC, ainsi qu'une evaluation des techniques d'indexation et d'accès utilisées.

### **V.1- Description du systeme d'accès aux donnees volumineuses**

Les traitements du systeme d'accès aux donnees volumineuses ont ete envisages comme un ensemble de modules effectuant, chacun une tâche spécifique.

Les principaux modules developpes sont les suivants:

- 1- Indexation
- 2- Acces aux donnees volumineuses
- 3- Mise a jour de la Base de Donnees Documentaires (BDD)

#### **v.1.1- Module d'indexation**

Une base de donnees documentaire est une collection de notices bibliographiques, chacune etant associee à un document (livre, these, article,...). Une notice contient toutes les informations qui permettent l'identification du document et de son contenu. Ces informations se presentent sous forme de texte de taille indéfinie et sont contenues dans differents champs, tels que le champ titre, auteur, resume, etc. Le nombre de champs descriptifs dans une notice n'est pas fixe, il varie d'un type de document a un autre.

Après avoir choisi la technique d'indexation appropriée, nous allons voir quelles sont les options a indexer. En effet, il existe **deux** approches pour indexer les informations

contenues dans une base de données documentaires ou referothèque:

- La première consiste à indexer chaque mot qui apparaît dans chaque champ d'une notice bibliographique et à constituer un lexique général unique. C'est le cas du système d'inversion totale présenté dans [GET 87].
- La seconde approche n'indexe que certains mots qui apparaissent dans certains champs et construit ainsi des lexiques spécialisés.

La première approche assure une recherche de haute précision, car elle permet de retrouver chaque mot paru dans la referoteque. Toutefois, cette méthode génère des index de mots de très grande taille [BLU 87].

La seconde approche limite l'indexation à certains types de mots de certains champs et permet une recherche sur la base des critères qui correspondent aux champs des mots qui sont indexés.

Comme le système d'accès aux données volumineuses est orienté vers une recherche de haute précision et comme la technique du Q-arbre est applicable pour des index de très grande taille, alors c'est la première approche qui a été retenue pour la conception et la réalisation du module d'indexation. Cependant, dans le cadre de notre projet et pour une première version du système, nous allons nous limiter aux champs titre, résumé et mot cle. En général, les résumés ont plus de mots clés que les titres. Par conséquent, si les mots clés sont sélectionnés à partir des résumés, la taille du fichier inverse devient importante et dépasse celle du fichier

de données lui-même. Et si les mots clés sont sélectionnés uniquement à partir des titres, alors, beaucoup de documents pertinents par rapport à la question et qui existent dans le fonds documentaire ne peuvent être retrouvés, ce qui donne un taux de silence élevé lors de la recherche [CHA 80] [DAL 89]. Par conséquent, la méthode la plus appropriée et réalisable qui permet d'extraire les mots clés à partir d'une grande masse de données, est celle de la suppression de mots vides.

Ce module se résume donc à :

1- Extraire tous les mots des champs titre, résumé et mot clé, et ceci, pour l'ensemble des notices du fichier bibliographique.

Dans cette étape nous indiquons aussi pour chaque mot :

- sa fréquence d'apparition dans la référentiel
- les informations de localisation du mot : le nombre et les adresses des notices qui le contiennent.

2- Filtrer les mots. Pour ce faire, nous appliquerons la notion de mots vides à la liste des mots obtenue à la première étape, et ceci en se référant à une liste de mots vides prédéfinie, établie à partir de l'ensemble des articles, prépositions, conjonctions, etc.

À la sortie de ce module, nous aurons une liste complète et ordonnée de tous les mots contenus dans les trois champs titre, résumé et mot clé et qui vont être réellement indexés, ainsi que les listes des notices qui les contiennent (une liste par mot). Il s'agit d'une **inversion totale**.

### v.1.2- Module d'accès aux données

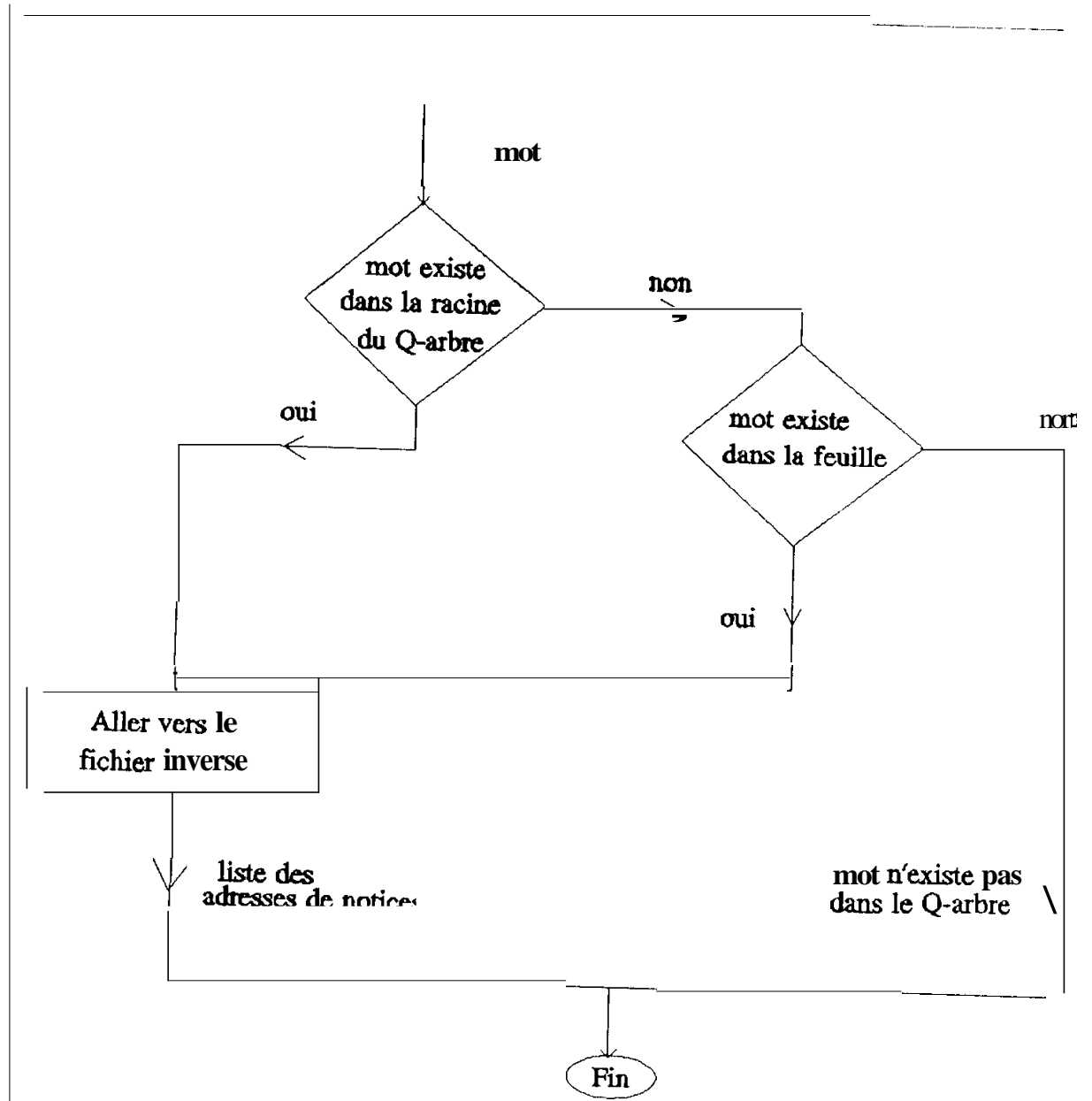
L'objectif de ce module est d'accéder aux données d'une référentielle en utilisant la technique du Q-arbre. Pour ce faire, il se divise en deux principaux sous-modules:

- Le premier sous-module consiste en l'implémentation du Q-arbre à partir de la liste des mots fournie par le module d'indexation. Pour cela, nous allons prendre parmi la liste des mots à indexer, ceux qui sont les plus fréquents et nous les mettons dans la racine du Q-arbre; ils seront donc en mémoire centrale et le reste des mots sera placé dans les feuilles du Q-arbre, c'est-à-dire en mémoire secondaire.

Le nombre de mots à mettre en mémoire centrale sera un paramètre à déterminer en faisant un compromis entre l'espace mémoire occupé et le temps de réponse du système. En fait, cela dépend essentiellement de la taille de la mémoire qui sera mise à la disposition du système. Toutefois, cette taille pourrait être augmentée en utilisant une méthode de compression au Q-arbre.

- Le deuxième sous-module établit une procédure de recherche dans le Q-arbre qui permet, pour un mot clé donné, de sélectionner les références bibliographiques contenant ce dernier.

Le schéma suivant est une illustration du chemin suivi dans le Q-arbre, pour l'exécution d'une commande de recherche:



**Schéma d'exécution d'une commande de recherche**

### **V.1.3- Module de mise a jour de la BDD**

La base annuelle d'informations est completee périodiquement par l'adjonction de la base documentaire mensuelle. Par consequent, le but de ce module est de mettre a jour mensuellement les divers fichiers rassemblant les informations sur lesquelles porteront les interrogations, a savoir, les fichiers inverse et lexicque associes a la base annuelle, ainsi que les structures de donnees associees a la technique d'accès, c'est-à-dire au Q-arbre.

## **V.2- Integration du systeme d'accès au systeme SIBADOC**

Etant donnée la structure modulaire du systeme documentaire SIBADOC (voir annexe a), l'intégration du systeme d'accès a ce dernier n'a pose aucun probleme d'implémentation.

Nous presentons dans ce qui suit la localisation dans le systeme SIBADOC des modules du systeme d'accès aux donnees volumineuses.

### **V.2.1- Module d'indexation**

Le module d'indexation est intégré au module de creation et mise a jour et il intervient au debut de cette operation. L'indexation agit sur le fichier bibliographique et permet d'extraire pour chaque notice l'ensemble des mots apparaissant dans les champs titre, resume et mot cle. L'indexation permet, ainsi, de creer une liste de couples, chacun etant constitué d'un mot cle et de l'adresse de la notice contenant ce mot cle.

Ce module consiste donc à créer pour chaque fichier bibliographique mensuel un fichier lexique contenant l'ensemble des mots extraits à partir des notices bibliographiques et un fichier inverse correspondant, contenant les adresses de ces notices.

#### **V.2.1.1- Structure du fichier lexique**

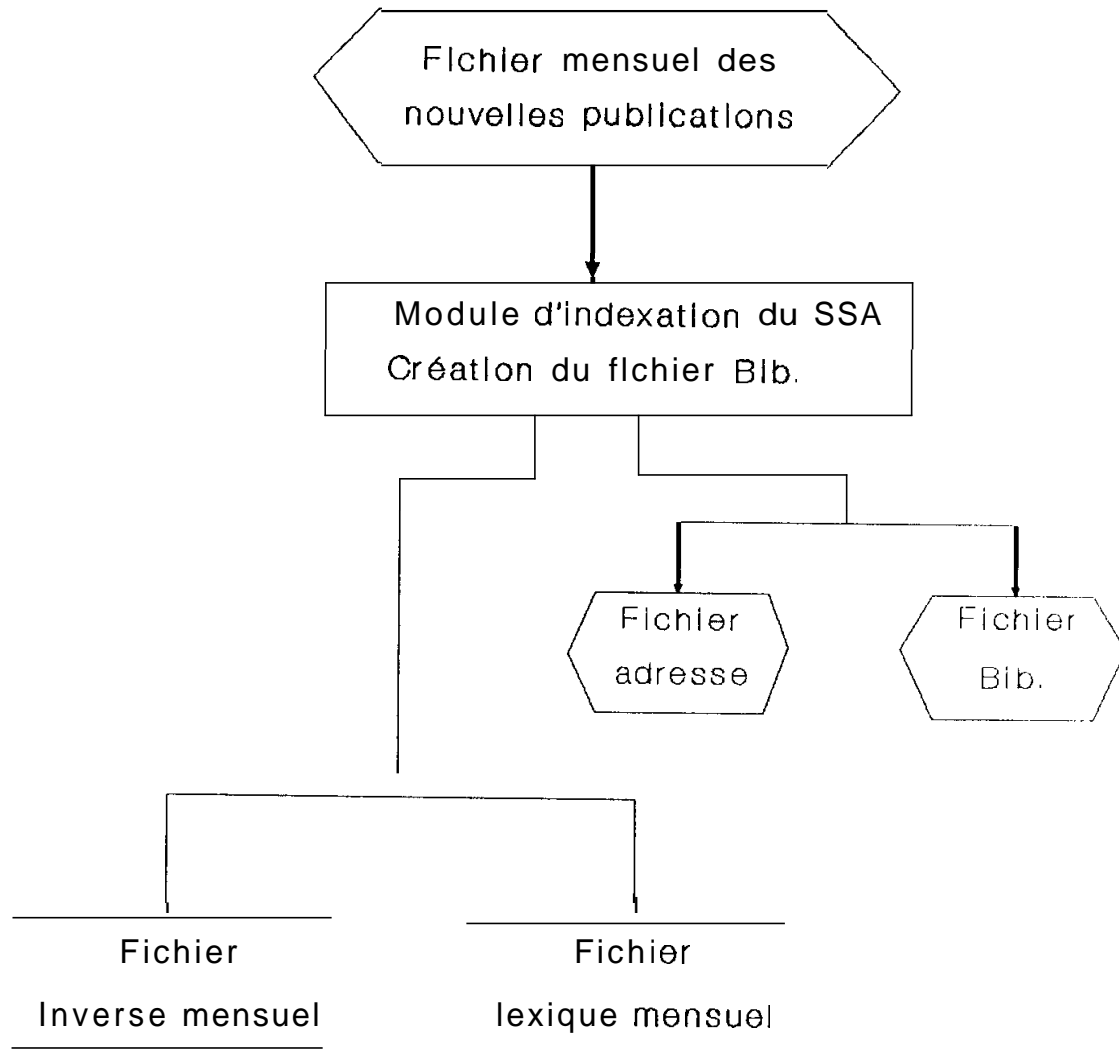
Le fichier lexique est un ensemble d'enregistrements présentes en ordre alphabétique et contenant chacun un mot cle et d'autres informations concernant ce dernier telles que sa longueur, sa fréquence d'apparition dans le corpus documentaire, un pointeur vers le fichier inverse, etc.

Une organisation en "séquentiel indexé" de ce fichier permet un accès rapide à l'information.

#### **V.2.1.2- Structure du fichier inverse**

Il contient pour chaque mot cle du fichier lexique une liste d'adresses des notices bibliographiques contenant ce mot clé.

Le schéma d'exécution du module de création et de mise à jour du système SIBADOC devient:



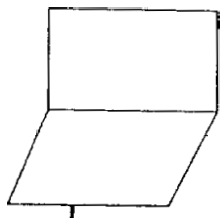
**Module de creation et mise a jour  
de la base avec indexation du sous systeme d'accès**

**V.2.2- Module d'accès aux données**

Ce module qui utilise la technique du Q-arbre pour l'exécution d'une commande de recherche, est intégré au module d'interrogation de la base et plus précisément au niveau de l'interpréteur de commandes.

Le processus d'interrogation n'ayant subi aucun changement du point de vue étapes de traitement, il se schématisera comme suit:

Utilisateur



Recherche utilisant  
la technique du Q-arbre

liste  
des adresses  
de notices

Module d'édition

notices Bib  
demandoos

Base documentaire

**Module d'interrogation de la base  
utilisant la technique d'accès du Q-arbre**

### **v.2.3- Module de mise a jour de la BDD**

Ce module est integre au module de creation et mise a jour du systeme **SIBADOC**. Il est exploitable en mode batch, ainsi que tous les autres modules du systemes d'accès aux donnees volumineuses, a l'exception du sous-module de recherche documentaire qui lui est executable en mode interactif.

### **v.3- Evaluation du systeme d'accès aux donnees volumineuses**

Après avoir integre le systeme d'accès aux donnees volumineuses au systeme **SIBADOC**, une evaluation de l'efficacité des techniques d'indexation et d'accès définies et utilisees, est presentee.

Cette evaluation se divise en deux grandes parties:

- Evaluation de l'indexation
- Evaluation de la technique d'accès

#### **v.3.1- Evaluation de l'indexation**

L'indexation est l'opération par laquelle on choisit les termes les plus appropries pour représenter le contenu d'un document. C'est donc l'opération centrale de tout systeme documentaire pour le stockage et la recherche des informations. Par consequent, l'indexation doit offrir une description pertinente des documents. Cette pertinence peut se mesurer essentiellement par l'exhaustivité et la specificite [PRA 69] [KIN 87] [LES 68].

L'exhaustivité caracterise la capacite a signaler toutes les notions importantes d'un document.

La spécificité caractérise la capacité à garder une information précise et non généralisée.

Plusieurs caractéristiques principales se dégagent de l'indexation utilisée par le système d'accès aux données volumineuses. Ces caractéristiques sont :

#### **- Evolutivité**

L'évolutivité permet de prendre en compte rapidement et rétrospectivement des évolutions des documents scientifiques et techniques. Ce concept est facilité par notre système, du fait que ce dernier s'appuie non pas sur un vocabulaire pré-établi, mais sur le vocabulaire extrait directement des documents.

#### **- Spécificité**

La méthode d'indexation utilisée par le système d'accès aux données volumineuses rend compte d'un niveau élevé de spécificité car elle rend compte de l'information telle qu'elle est citée dans les textes, sans effet de généralisation qui est lié à un effet de synthétisation de l'information. En effet, certains termes généraux telles que intelligence artificielle, reconnaissance des formes, etc, sont introduits de façon quasi systématique par les indexeurs dans le but de traduire le contenu des documents pour des demandeurs potentiels travaillant dans des disciplines variées, ce qui n'est pas le cas de notre système.

#### **- Exhaustivité**

L'indexation utilisée par le système d'accès aux données volumineuses présente une bonne exhaustivité, du fait qu'elle

prend en compte tous les mots du corpus documentaire (du moins important au **plus** important).

Outre ces caractéristiques, deux principales mesures sont utilisées pour évaluer l'efficacité d'une recherche documentaire, et qui sont le taux de rappel et le taux de précision.

#### ▪ **Definition du taux de rappel**

Le taux de rappel est la proportion des documents pertinents retrouvés par rapport à l'ensemble des documents pertinents présents dans le système. Par exemple, si, dans une base de données de 100 références, 20 répondent à une question, et que le système en retrouve 15 mais omet les 5 autres, le taux de rappel est 15/20 ou 75%. Les documents pertinents qui n'ont pas été retrouvés constituent ce que l'on appelle le silence.

$$\text{Taux de rappel} = \frac{\text{Nombre de documents pertinents retrouvés}}{\text{Nombre de documents pertinents existants}}$$

#### ▪ **Definition du taux de précision**

Le taux de précision ou de pertinence est la proportion des documents pertinents par rapport à l'ensemble des documents fournis par le système. Si une recherche donne 40 références en réponse et que 15 documents sont pertinents,

le taux de precision est 15/40, c'est-à-dire 37,5%. On appelle bruit les 25 documents non pertinents qui ont ete fournis en même temps.

$$\text{Taux de precision} = \frac{\text{Nombre de documents pertinents retrouves}}{\text{Nombre de documents extraits}}$$

Rappel et precision sont des qualités contradictoires. Plus la recherche est precise, plus le risque est grand qu'elle ne retienne pas les documents qui auront ete decrits en terme plus generaux mais qui peuvent concerner le sujet, donc le rappel diminue et le silence augmente. Inversement, plus la recherche est large, de façon a retrouver le plus grand nombre possible de documents ayant un rapport, si limité soit-il, avec le sujet, plus le risque est grand de trouver en même temps des documents non pertinents, donc d'augmenter l'imprecision et le bruit.

La solution est d'avoir une moyenne raisonnable de ces deux parametres. C'est-à-dire ni le taux de rappel, ni le taux de precision ne doivent être faible en même temps et donc avoir un vocabulaire d'indexation ni trop general ni trop spécifique. En pratique une indexation ayant 60% de taux de rappel et 50% de taux de precision est satisfaisante pour les utilisateurs.

Comme notre systeme est base sur une inversion totale, l'interrogation est tres souple et peut être menée de façon a privilégier soit le taux de rappel, soit le taux de precision [SWE 63]. En effet, l'équation de recherche conduit a un mauvais taux de rappel si elle ne réussit pas a

couvrir les différents aspects possibles de la question, utilise une formulation **trop spécifique** ou trop exhaustive (on ne trouvera pas assez de documents répondant à tous ces critères). Elle conduit à une mauvaise précision si elle n'est pas assez spécifique et exhaustive, si elle utilise des termes ou des combinaisons de termes inadéquats et si sa logique est defectueuse.

### **v.3.2- Evaluation de la méthode d'accès**

Dans ce paragraphe une comparaison entre la technique d'accès utilisée par le système SIBADOC et celle du Q-arbre sera présentée. Cette comparaison prend en compte les paramètres temps CPU et nombre d'entrée/sortie pour la recherche.

Pour la comparaison nous avons utilisé un échantillon de 3000 notices de la base de données INSPEC de l'année 1990, soit une liste de 18996 mots clés.

Pour rechercher un mot clé, le système SIBADOC qui utilise une technique d'inversion classique doit accéder au fichier répertoire, au fichier lexique et au fichier inverse, ce qui nécessite au minimum trois accès disque. Par contre pour la technique du Q-arbre, ou bien le mot clé existe dans la racine et donc il ne nécessitera aucun accès disque car la racine est en mémoire centrale, ou bien il est dans une des feuilles du Q-arbre et il nécessite dans ce cas un seul accès disque pour le retrouver. Par conséquent, la technique du Q-arbre nécessite au plus un accès disque pour rechercher un mot clé donné.

Le temps de recherche est primordial, pour les utilisateurs, dans un système conversationnel. Les résultats expérimentaux montrent que la technique du Q-arbre est plus rapide que celle de l'inversion utilisée par le système **SIBADOC**, et ceci est une conséquence du nombre d'accès disque utilisé par chacune d'elles lors de la recherche.

De plus, lorsque le nombre de mots dans la racine augmente le temps de recherche diminue (voir tableau ci-dessous des résultats). Aussi, pour augmenter les performances de la technique du Q-arbre il faut utiliser une technique de compression pour la racine, ce qui augmente le nombre de mots à mettre dans cette dernière.

Quelques résultats de la comparaison sont illustrés par le tableau suivant:

Mots clés	SIBADOC temps cpu ( $10^{-2}$ s)		
		500 mots dans la racine	1500 mots dans la racine
Robot	25	21	13
Amos	22	18	11
Pransient	25	21	11
Unix	23	21	11
Apur	22	19	00
Vacuum	25	22	09
Vibration	25	23	12
Voltage	24	22	11
Workstation	25	23	11
Wavelength	24	22	10

## *CONCLUSIONS ET PERSPECTIVES*

L'optimisation du stockage et de l'accès aux données volumineuses telle qu'elle est décrite dans ce document est d'un grand secours pour l'exploitation d'un grand volume d'informations et pour le repérage rapide de documents, ce qui est particulièrement important.

Le système d'accès aux données volumineuses, dont la conception et la réalisation ont constitué l'essentiel de notre travail, est organisé de façon modulaire. Il est adaptable à tous les systèmes gérant de grands volumes de données textuelles et est basé sur des techniques d'indexation et d'accès qui ont les caractéristiques suivantes :

- Facilité de mise en œuvre. En effet, la technique du Q-arbre et la méthode d'indexation utilisées sont très simples à mettre en œuvre.
- Le vocabulaire d'indexation est structuré à partir des documents, il évolue donc en même temps que le fonds documentaire. De plus, comme l'indexation est basée sur une inversion totale, prenant en compte tous les mots du corpus documentaire, l'interrogation est très souple et peut être menée de façon à privilégier soit le taux de rappel, soit le taux de précision.
- La technique d'accès étant basée sur le Q-arbre, elle est très rapide du fait qu'elle nécessite, au plus, un accès disque pour rechercher le mot désiré. De plus, elle est applicable à des index de très grande taille, ce qui a permis l'utilisation de la technique d'inversion totale et donc l'augmentation de l'efficacité de la recherche d'information.

Toutefois, nous ne pouvons pas prétendre avoir traité tous les problèmes liés à l'accès aux données volumineuses. En effet, plusieurs extensions à notre travail peuvent être envisagées :

- Prendre en compte les mots clés composés : un système qui n'utilise pas les mots composés est caractérisé par une grande souplesse due à la simplicité de mise en œuvre et à la facilité des mises à jour. Il est caractérisé aussi par une indexation fine, mais présente une certaine perte dans la précision.

La prise en charge des mots composés peut être réalisée de deux façons :

- Soit en utilisant la notion de mots composés dans le lexique. Par exemple, considérer "Informatique documentaire" comme un seul mot.

- Soit en utilisant, en plus des opérateurs classiques ET/OU, les opérateurs de voisinage tel que l'opérateur ADJACENT qui est utilisé pour préciser la position relative des termes entre eux.

- Utiliser une méthode d'indexation automatique se basant sur l'approche linguistique et où tous les problèmes liés à l'indexation seront pris en compte (synonymie, polysemie, normalisation des termes, pertinence des termes, relations entre termes d'indexation, ...).

- Appliquer une méthode de compression au Q-arbre afin d'augmenter le nombre de mots à mettre en mémoire centrale, et par conséquent augmenter les performances de la méthode d'accès.

- Offrir la possibilité de rechercher un mot en fonction de sa position: par exemple, dans le titre, dans le resume, dans tel ou tel paragraphe, etc.

- Prendre en compte les donnees ecrites dans des langues autres que l'Anglais, par exemple, le Francais ou l'Arabe. Ceci est tres simple a realiser du fait de la modularité du systeme. En effet, il suffit d'avoir un dictionnaire de mots vides de la langue desiree (Francais ou Arabe) est faire de légère modification au niveau du module d'indexation et plus precisement au niveau de l'extraction des mots cles.

- Intégrer le systeme d'accès aux donnees volumineuses au systeme de gestion d'objets complexes SYGOC. En effet, les techniques definies et utilisées dans notre systeme d'accès sont facilement adaptables a un environnement autre que celui des bases de donnees documentaires: c'est, par exemple, le cas de l'environnement des objets complexes (formulaires, images, diagrammes, etc), qui fait a l'heure actuelle l'objet de nombreux travaux de recherche.

*BIBLIOGRAPHIE*

- [AHO 75] : AHO A.V. and CORASICK M.J.  
Efficient string matching.  
Commun. ACM 18, p. 333-340, 6 (June 1975).
- [ALS 75] : ALSBERG, A.PETER  
Space and time savings through large database  
compression and dynamic restructuring.  
Proc. of the IEEE, vol.83, N°8, Aug 1975,  
pp 114-1122.
- [AST 76] : ASTRAHAN M.M., and CHAMBERLIN D.D.  
SYSTEM R: Relational approach to database  
management  
ACM TRANSACTIONS OF DATABASE SYSTEMS
- [BAN 88] : BANCILHON F.  
Object oriented database systems.  
Int. Symp. on principles of database systems,  
Austin (Texas), 1988.
- [BAY 72] : BAYER R., and McCREIGHT E.  
Organization and maintenance of large ordered  
indexes.  
ACTA informatica vol.1, n°.3, p. 173-189, 1972.
- [BEL 87] : BELKIN N.J.  
Retrieval techniques  
Annual Review of Information Science and  
Technology (ARIST), vol.22, 1987.

- [BEN 88] : BENDAKIR M., and YAHIA AISSA N.  
Description du systeme documentaire SIBADOC.  
CERIST, 1988.
- [BEN 89] : BENDAKIR M., and YAHIA AISSA N.  
Le systeme **SIBADOC**.  
CERIST, 1989.
- [BLA 85] : BLAIR D.C., and MARON M.E.  
An evaluation of retrieval effectiveness for a  
full-text document-retrieval system.  
Communication of the ACM, **vol.** 28, n°.3,  
p 289-299, March 1985.
- [BLU 87] : BLUMER A., and HAUSSLER D.  
Complete inverted files for efficient text  
retrieval and analysis.  
ACM vol. 34, n°.3, **p.** 578-595, Jul. 1987.
- [BOU 91] : BOUMALHA H., BENDJENAH F.  
Vers la réalisation d'un systeme de gestion  
d'objets complexes: SYGOC  
Dans ler Séminaire sur les Bases de Donnees,  
p 216-227, Juin 1991, Alger
- [BOY 77] : BOYER R.S., and MOORE J.S.  
A fast string searching algorithm.  
Commun. ACM 20, P. 762-772, 10 (**Oct.** 1977).

- [CHA 73] :       **CHAUMIER J.**  
Systemes informatiques de documentation.  
Entreprise moderne d'edition, 1973.
- [CHA 80] :       **CHAUMIER J.**  
Travail et methode du/de la documentaliste.  
Collection formation permanente en  
sciences humaines. 1980.
- [CHA 90] :1      **CHAUMIER J., and DEJEAN M.**  
L'indexation documentaire.  
De l'analyse conceptuelle humaine a l'analyse  
automatique morphosyntaxique.  
Documentaliste, vol.27, n°6, Nov.- Dec. 1990.
- [CHR 84] :       **CHRISTODOULAKIS S. , and FALOUTSOS C.**  
Design considerations for a message file server.  
IEEE Trans. Softw. Eng. SE- 10, p.201-210,  
2 (Mar. 1984)
- [CHI 86] :       **CHIARAMELLA Y.**  
**IOTA: A** full text information retrieval system.  
Organisation of the **1986-ACM** Conference on  
Research and Development in information. 1986.
- [COD 70] :       **E. F. CODD**  
The relational model of data for large  
shared data banks.  
**Comm.** of **ACM**, vol.13, n°6, Juin 1970.

- [COD 79] : E.F.CODD  
Extending the database relational model  
to capture more meaning.  
IBM research laboratory.
- [COO 78] : COOPER W.S. and MARON M.E.  
Foundations of probabilistic and utility  
theoretic indexing.  
Journal of the Association for Computing  
Machinery, vol.25, n°1, p.67-80, Jan. 1978.
- [COY 67] : COYAND Maurice, SIOT-DECAUVILLE Nelly  
L'analyse automatique des documents.  
Publié par le Centre National de la  
Recherche Scientifique, Paris, 1967.
- [DAL 89] : DALBIN S. , and CHARTRON G.  
Indexation manuelle et indexation automatique.  
Documentaliste vol. 26, n°.4-5, p.181-187  
(Jui. 1989).
- [DEW 81] : DEWEZE, A.  
Reseaux sémantiques : Essai de modelisation -  
Application a l'indexation et a la recherche de  
l'information documentaire.  
These de doctorat, Lyon 1981.
- [DEW 83] : DEWEZE, A.  
L'accès en ligne aux bases documentaires.  
Collection MASSON, 1983.

- [DEW 85] : DEWEZE, A.  
Informatique documentaire.  
Collection MASSON, 1985.
- [DEW 89] : DEWEZE A.  
Information documentaire. Methode + Programme.  
Collection MASSON, 1989,
- [EAS 89] : EASTMAN C.M.  
Approximate retrieval : A comparison of  
information retrieval and database management  
system in database engineering, vol.12, n°.2,  
p.41-45, 1989.
- [FAL 84] : CHRIS FAMUTSOS and STAVROS CHRISTODOULAKIS  
Signature files: an access method for documents  
and its analytical performance evaluation.  
ACM Transactions on Office Information Systems,  
vol. 2, n°.4, p. 267-288, Octo. 1984.
- [FAL 85] : CHRIS FALOUTSOS  
Access methods for text  
Computing Surveys, vol.17, n°.1, March 1985
- [FUM 82] : FUMIRO.MATSMO and SHOUICHI.FUTAMURA  
A Quasi-Optimum text compression code.  
Ibid., vol.55, N°2, 1982, pp 103-106.

- [FUM 86] : FUMIHIRO MATSUO, SHOUICHI FUTAMURA, and TAKESHI SHINOHARA.  
Efficient storage and retrieval of large document database.  
computer, **p.** 456-463, jun. 1986.
- [GAR 85] : GARDARIN G. , and VALDURIEZ P.  
Bases de donnees relationnelles  
Analyse et comparaison des systemes. 1985
- [GAR 89] : GARDARIN **G.**, and VALDURIEZ **P.**  
SGBD relationnels: Analyse et comparaison de bases de donnees. 1989
- [GAR 91] : **GARDARIN G.** , and **VALDURIEZ P.**  
**SGBD avances:**  
Bases de donnees objets, deductives, réparties.  
**1991**
- [GET 87] : GETHIN P.  
Text retrieval.  
In : Database and Network Journal, vol.17, n 3,  
**p. 2-7, 1987.**
- [GOT 75] : GOTTLIEB, DORON, HAGETH, A.STEVEN, LEHOT,  
G.H.PHILIFE **and** RABIBOWITZ, S.HENRY  
A classification of compression method and their usefulness for large data processing center.  
National Computer Conference 1975, pp 453-458.

- [GOY 85] : GOYEL P., and KOTAMARTI V.  
A strategy for compressed storage and retrieval  
of documents.  
Compressed storage and retrieval.  
p. 224-232, 1985.
- [GUI 89] : GUINCHAT Clair, SKOURI Yolande  
Guide pratique des techniques documentaires.  
Volume 2. Traitement de l'information ".  
Diffusion EDICEF, 1989.
- [GUI 90] : GUINCHAT Clair, NENOV Michel  
Sciences et techniques de l'information et de  
la documentation.  
Publié par l'Organisation des Nations Unies pour  
l'Education, la Science et la Culture (UNESCO),  
1990.
- [HEA 72] : **HEAPS. H.S**  
**Storage** analysis for a compression coding for a  
document databases.  
INFOR, VOL.10, N°1, Feb 1972, pp 47-61.
- [HOL 79] : HOLLAR L.A.  
Text retrieval computers.  
IEEE computer 12, p. 40-50, 3 (Mar. 1979).
- [IAE 82] : IAEA-INIS-17  
First step on STAIRS.  
Document interne, 1982.

- [IBM 79] : IBM WORLD **TRADE CORPORATION**.  
IBM System/370 (OS/VS).  
Storage and Information Retrieval  
System/Vertical Storage.  
Reference Manuel.
- [JOH 69] : JOHN R. , and *HARRY D. HUSKEY*  
An information retrieval system based on  
superimposed coding.  
In Fall Joint Computer Conference, **1969**.
- [KEV 83] : KEVIN P.J.  
How do we index?: A report of some ASLIB  
informatics group activity.  
Journal of documentation vol. **39**, n°.1,  
Mar.**1983**.
- [KIM 87] : KIM W., and **BANERJEE J.**  
Data model issues for object oriented  
applications.  
ACM TOIS, vol.5, n°1, 1987.
- [KIN 87] : KINNUCAN M. and NELSON M.  
Statistical methods in information science  
research.  
Annual Review of Information Science and  
Technology (**ARIST**) , vol.22, **1987**.
- [KNU 77] : **KNUTH D.E.**, **MORRIS J.H.** and **PRATT V.R.**  
Fast pattern matching in strings.  
SIAM J. Comput. 6, p. **323-350**, 2 (June **1977**).

- [KUO 82] : KUO-CHUNG TAI, and ALAN L. THARP  
The practicality of text signatures for  
accelerating string searching.  
Software-Practice and Experience, vol. 12,  
p.35-44, 1982.
- {LAL 87} : LALLICH-BOIDIN, ROUAULT, J.  
Un modele linguistique de reconnaissance du  
français pour un systeme documentaire.  
CRISS Université des Sciences Sociales de  
Grenoble, 1987.
- [LAM 82] : LAM X. and SALTON G.  
Term weighting in information retrieval using  
the term precision model.  
Journal of the Association for Computing  
Machinery, vol.29, n°1, p.152-170, Jan. 1982.
- [LEC 88]: LECLUSE C.  
O2 an object oriented data model.  
Proceeding of the ACM SIGMOD conference,  
Chicago 1988,
- [LES 68] : LESK S.  
Computer evaluation of indexing and text  
processing.  
Journal of the Association for Computing  
Machinery, vol.15, n°1, Jan. 1968.
- [LON 80] : LONG B.  
Linguistique et indexation  
Documentaliste, vol. 17, n° 3, Mai-Juin 1980

- [MAL 71] : **MALCOLM C. HARRISON**  
Implementation of the substring test by hashing.  
Communications of the **ACM** vol. 14, n°.12,  
p. 777-779, Dec. 1971.
- [MAR 90] : **MARON M.E., and BLAIR D.C.**  
Full-text information retrieval: further  
analysis and clarification.  
Information Processing & Management vol. 26,  
n°.3, p. 437-447, 1990.
- [MAT 81] : **MATSUO F., FUTAMURA S., TAKAGI T., YOSHIDA S.**  
Organization of word dictionary indexes for  
quick retrieval.  
Ibid., vol. 54, n°.3, p. 183-187 1981.
- [MCG 83] : **MCGILL M.J., and SALTON G.**  
Introduction to modern information retrieval.  
McGRAW-HILL, **New York, 1983.**
- [MEC 89] : **MECHKOUR M.**  
Integration d'une technique de compression de  
donnees a un système de recherche documentaire.  
Mémoire de fin d'études, INI **1989.**
- [MER 84] : **MERRETT T.H.**  
First steps to algebraic processing of text.  
**New Applications of Databases.**  
**ISBN 0-12- 275550-2, 1984.**

- [MET 88] : METZGER Jean-Paul  
Syntagmes nominaux et information textuelle.  
Reconnaissance automatique et representation ".  
These de docteur d'état en sciences.  
Universite Claude Bernard, Lyon1, 1988.
- [PEC 85] : PECHURA, A.MICHAEL and DAVID R.MCINTYRE  
Data compression using static HUFFMAN  
code-decode tables.  
Communications of the ACM, vol.28, N°6,  
June 1985, pp 612-616.
- [PIK 81] : PIKE J.  
Text compression using a 4-bit ending schema.  
The computer journal, vol 24, N°4, 1981,  
pp 324-330.
- [PRA 69] : PRANAS ZUNDE and MARGARET E. DEXTER  
Indexing consistency and quality  
American Documentation. July 1969.
- [RAB 84] : RABATTI F., and ZIZKA J.  
Evolution of access methods to text documents in  
office systems.  
In Proceedings of the 3rd Joint ACM-BCS  
Symposium on Reseach and Development in  
Information Retrieval. 1984.
- [REG 81] : REGHBATI, HASSEN K.  
An overview of data compression techniques.  
Computer Vol.14, N°4, April 1981, pp 71-75.

- [RIV 90] : RIVIER A.  
Construction des langages d'indexation,  
Aspects theoriques.  
Documentaliste, vol. 27, n° 6, Nov-Dec 1990
- [ROB 79] : CHARLES S. ROBERTS  
Partial match retrieval via the method of  
superimposed codes.  
In Proceedings of the IEEE, vol. 67, n°.12,  
Dec, 1979.
- [ROU 87] : ROUAULT Jacque  
Linguistique automatique. Application  
documentaires.  
Edition Peter Lang SA, Berne, 1987.
- [SAC 87] : SACKS-DAVIS R., and KENT A.  
Multikey access methods based on superimposed  
coding techniques.  
ACM Trans. on Database Sys., vol. 12, n°.4,  
p 655-696, Dec.1987.
- [SAL 70] : SALTON G.  
Automatic text analysis: Automatic document  
indexing and classification methods are examined  
and their effectiveness is assessed.  
Science, vol.168, April 1970.

- [SAL 71] : SALTON G.  
The SMART Retrieval System-Experiment in  
automatic document processing.  
Prentice-Hall, Englewood Cliffs, New Jersey.  
1371
- [SAL 86a] : SALTON G.  
Recent trends in automatic information  
retrieval.  
Proceeding of the 1986-ACM Conference on  
Research and Development in Information  
Retrieval. Pages 1-10. Pisa, Sept. 1986.
- [SAL 86b] : SALTON G.  
Another look at automatic text-retrieval  
systems.  
Communication of the ACM, vol. 29, n°.7,  
p 648-656, July 1986.
- [SAL 89] : SALTON G.  
Automatic text processing.  
The transformation, analysis, and retrieval  
of information by computer.  
Addison-Wesley publishing company, 1989.
- [STO 76] : STONEBRACKER, M, and WONG, E.  
The design and implementation of INGRES  
ACM Transactions on Database Systems, Sept. 1976

- [SWE 63]: SWETS John  
Information Retrieval System. Statistical  
decision theory may provide a measure of  
effectiveness better than measures proposed  
to date.  
Massachusetts Institute of Technology,  
Cambridge, 16 July 1963.
- [TAV 92] : TAVAKOLI N., and RAY A.  
A new signature approach for retrieval of  
documents from free-text databases.  
Information Processing and Management  
**vol.28**, n°2, p.153-163, 1992.
- [TIG 92] : TIGRINE B.  
L'Indexation  
Document interne au CERIST, 1992.
- [TSI 82] : TSICHRITZIS D. CHRISTODOULAKIS  
Message files., p. 110-112, **ACM** 1982.
- [VAL 87] : VALDURIEZ P.  
Objets complexes dans les systèmes de bases  
de données relationnels.  
TSI vol.6, n°5, **1987**.
- [WIS 86] : WISNIEWSKI, JANUSZ L.  
Compression of index dictionary in an inverted-  
file-oriented-database : Some effective  
algorithms.  
Info.Proceeding & management, vol.22, N°6, 1986.

[WIS 87] : WISNIEWSKI, JANUSZ L.  
Effective text compression with simultaneous  
bigram and trigram encoding.  
North Holland, Journal of information  
science 13, 1987.

.

*ANNEXES*

## A N N E X E A

SIBADOC : Un systeme d'interrogation des bases de donnees  
documentaires

- 1- Description du systeme SIBADOC
  - 1.1- Module de creation et de mise a jour
    - 1.1.1- Structure des bases bibliographiques
    - 1.1.2- Creation du lexique
    - 1.1.3- Inversion
  - 1.2- Module d'interrogation des bases documentaires
    - 1.2.1- Langage d'interrogation
    - 1.2.2- L'interpreteur des commandes

Une base de données documentaire est une collection de notices bibliographiques, chacune étant associée à un document (livre, thèse, article, ...). Une notice contient toutes les informations qui permettent l'identification du document et de son contenu. Ces informations se présentent sous forme de texte de taille indéfinie et sont contenues dans différents champs, tels que le champ titre, auteur, résumé, etc. Le nombre de champs descriptifs dans une notice n'est pas fixe, il varie d'un type de document à un autre.

Les systèmes documentaires ont pour objet, comme tous les systèmes de gestion de données, de répondre à un besoin d'information. Leur objectif n'est pas de transmettre l'information primaire, mais plutôt des informations décrivant plus ou moins bien cette dernière, enregistrées sous forme de notices bibliographiques.

Les systèmes documentaires sont donc des systèmes d'information dont l'objectif est de mettre à la disposition des utilisateurs les documents secondaires; ce qui leur permet dans un premier temps de procéder à une sélection sans devoir consulter les documents physiques, ensuite de localiser ces derniers.

Le système SIBADOC, développé au CERIST!, s'inscrit dans cette optique et permet la gestion et l'accès à l'information documentaire des bases bibliographiques internationales acquises par le centre auprès d'agences spécialisées. Actuellement, SIBADOC permet la gestion et l'accès aux bases bibliographiques suivantes:

INSPEC: pour la technologie et les sciences de l'ingénieur  
(INformation Services for the Physics and Engineering  
Communities).

INIS : pour le domaine du nucleaire et les techniques qui s'y  
rapportent (International Nuclear Information  
System)

AGRIS : pour le domaine de l'agriculture  
(AGRicultural Information System)

Ces bases sont commercialisees sous le format d'echange  
international ISO 2709.

Notons que toute base possedant une structure d'organisation  
similaire a celle des trois bases citees ci-dessus, est  
directement exploitable par le systeme SIBADOC.

## 1- Description du systeme SIBADOC

le systeme SIBADOC est constitue de deux principaux  
modules: le module de creation et de mise a jour des bases et  
le modules d'interrogation [BEN 88].

### 1.1- Module de creation et de mise a jour

#### 1.1.1- Structuration des bases bibliographiques

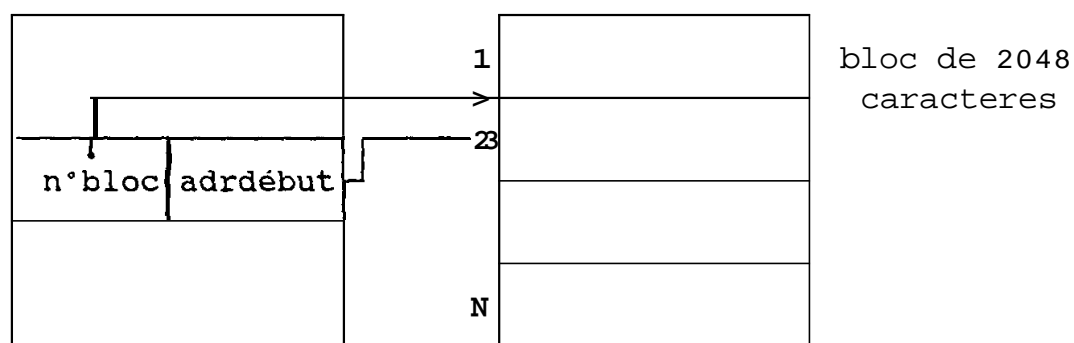
Cette premiere etape consiste a traiter les donnees  
emanant de differents producteurs, et donne en sortie un couple  
de fichiers par mois, en l'occurrence un fichier d'adresses et  
le fichier bibliographique.

## a- Structure du fichier d'adresses

L'enregistrement de ce fichier est de la forme (adr,depl) où "adr" désigne l'adresse du bloc et "depl" le déplacement ou l'adresse de la notice dans le bloc. L'objectif du fichier adresse est de localiser la notice dans le fichier bibliographique.

## b- Structure du fichier bibliographique

C'est un ensemble d'enregistrements de longueur variable (appelés notices), chacun étant repéré par une adresse double (numero de bloc, déplacement dans le bloc) dans le fichier adresse. Cet ensemble de notices est organisé par périodes. En effet, les notices du même mois sont regroupées dans le même fichier et les fichiers mensuels sont regroupés en bases annuelles dont la concatenation constitue la base générale.



Structure du fichier bibliographique mensuel

### 1.1.2- Creation du lexique

Pour les bases **INIS** et **AGRIS**, les thesauri existent et sont repts sur bandes magnetiques, par contre pour la base INSPEC, le CERIST dispose uniquement du thesaurus manuel.

Dans le cadre du projet **SIBADOC**, il a été procédé a la creation des lexiques pour les differentes bases.

#### 1.1.2.1- Le fichier lexique

Le fichier lexique est constitué d'un répertoire et d'un dictionnaire de mots ou expressions de la langue naturelle representant des descripteurs, appeles aussi mots cles et tries par ordre alphabétique. Ces descripteurs sont extraits a partir du champ "mot **clé**" du fichier bibliographique et constituent un ensemble de partitions dont chacune rassemble les mots qui ont une même initiale.

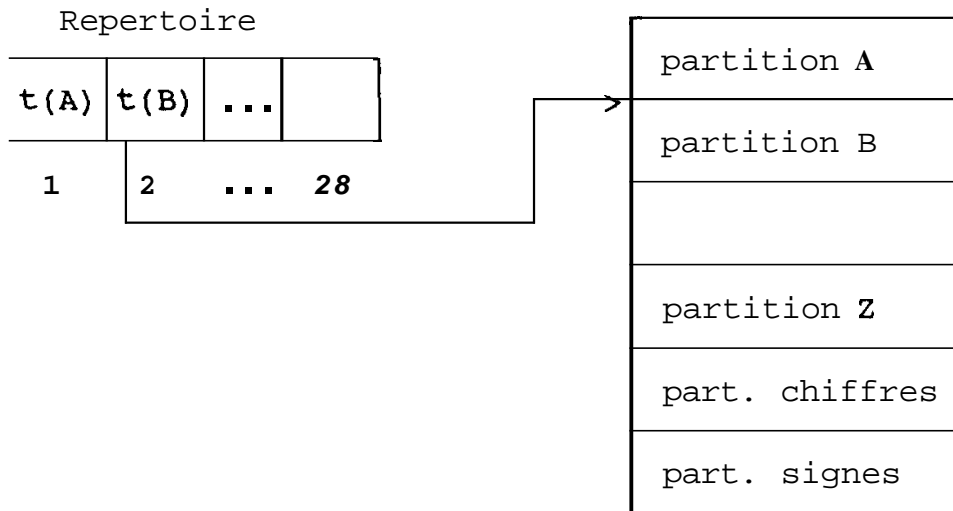
Il existe **28** partitions:

- **26** pour l'alphabet (A..Z)
- 1 partition por les chiffres (0..9)
- 1 partition pour les signes speciaux (-,\* , ...)

Le repertoire facilite la recherche dans une partition puisqu'il indique le debut et la taille de chaque partition.

Remaraue :

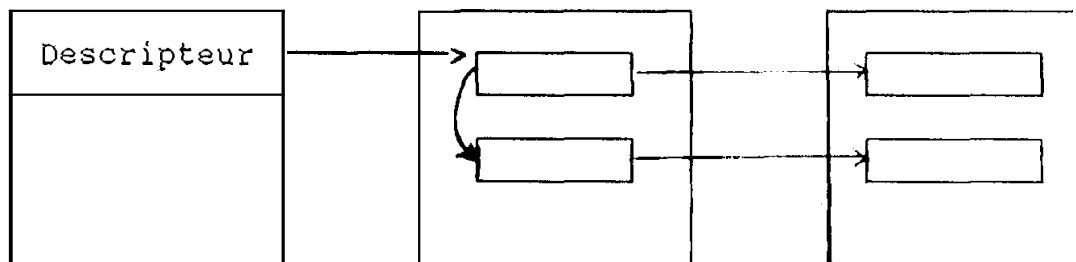
Nous avons un fichier lexique par an.



Structure du fichier lexique

1.1.3- Inversion

Cette phase consiste a creer le fichier inverse qui permet de faire le lien entre les descripteurs et les notices qu'ils indexent dans le fichier bibliographique.



Structure du fichier inverse

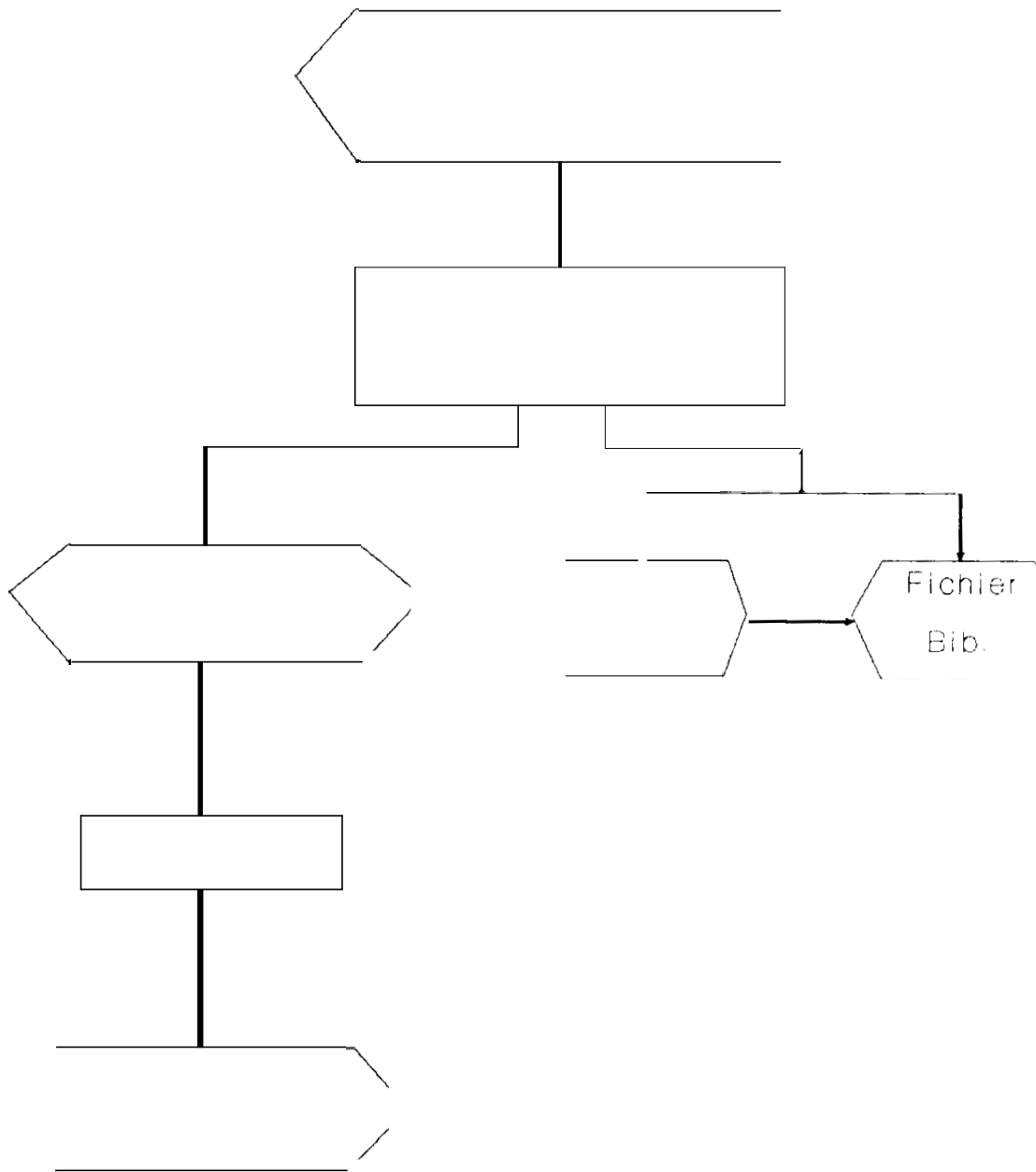
Remarque :

Nous avons un fichier inverse par an.

Le module de creation et de mise a jour permet donc l'évolution de la base par le traitement et l'insertion dans celle-ci du fichier mensuel des nouvelles publications. Pour cela, quatre operations sont effectuees :

- 1- Creation du fichier bibliographique et de son fichier adresse pour le mois courant.
- 2- Indexation des documents par les mots cles figurant dans les notices par la creation d'une liste de couples (mot clé, numero de notice contenant ce mot clé).
- 3- Inversion : creation du fichier inverse mensuel
- 4- Mise a jour des fichiers inverse et lexique annuels (uniquement l'ajout est pris en compte).

Le module de creation et de mise a jour *se* schématise comme suit :



**Module de creation et mise a jour de la base**

## 1.2- Module d'interrogation des bases documentaires

Ce module permet la consultation du fonds documentaire, et comporte principalement :

- Un langage d'interrogation
- Un interpreteur de commandes.

### 1.2.1- Langage d'interrogation

Ce langage offre de riches possibilités de recherche aux utilisateurs desireux de communiquer avec n'importe quelle base documentaire disponible au CERIST.

Il permet l'écriture d'une equation de recherche ou requête et de traiter tous les rapprochements informatiques nécessaires a sa satisfaction. Il offre également des services d'assistance, d'édition et de navigation dans la base (changement de bases et d'années de la base).

Ces differentes fonctions sont les suivantes [BEN 89] :

#### 1- Fonctions d'assistance

HELP : demande d'aide, fournit une documentation  
détaillee sur l'ensemble des commandes.

#### 2- Fonctions de recherche

FIND : fonction de recherche par mot cle

INDEX : fonction d'extension de racine, recherche tous  
les mots cles contenant la racine specifiée.

SELECT: fonction de selection d'un mot cle, d'une liste  
de mots clés, apparaissant apres l'exécution de  
la fonction index.

3- Fonctions d'édition/impression

DISPLAY : affichage des notices bibliographiques  
PRINT : impression des notices bibliographiques

4- Fonctions de service

CHANGE : permet de changer de base ou d'année dans  
la base choisie et de charger les fichiers  
correspondants.  
QUIT : terminaison de la session

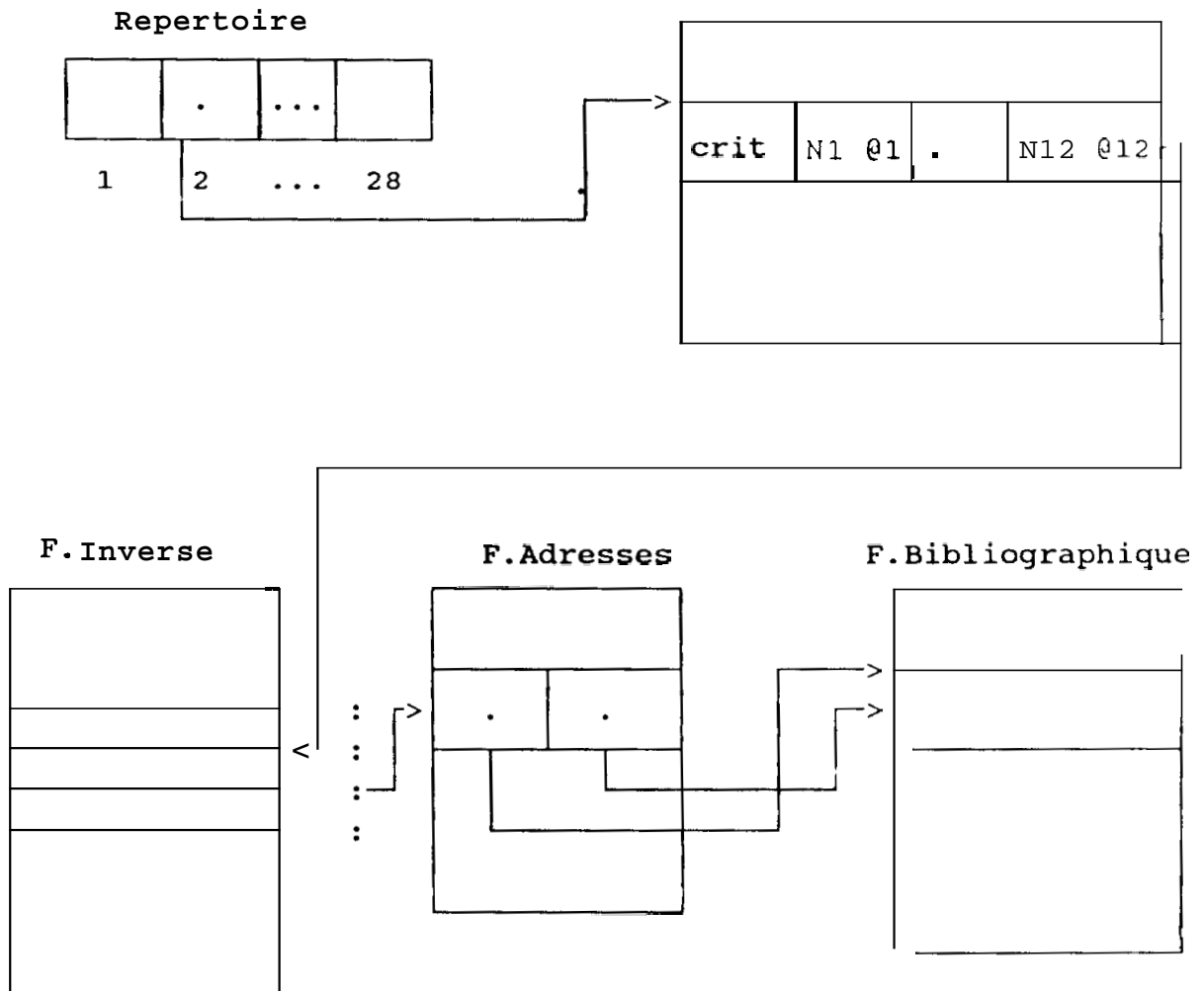
6- Operateurs

OR, AND, NOT, \$ (pour l'extension) et #(pour le masquage).

Une requête est un ensemble de descripteurs liés entre eux par des relations logico-syntaxiques issues de la syntaxe du langage documentaire propre à SIBADOC. Par conséquent, toute requête établie par un utilisateur doit être soumise à une analyse syntaxico-sémantique à l'issue de laquelle celle-ci sera codée dans une notation interne et transmise à l'interpréteur de commandes.

1.2.2- L'interpréteur de commandes

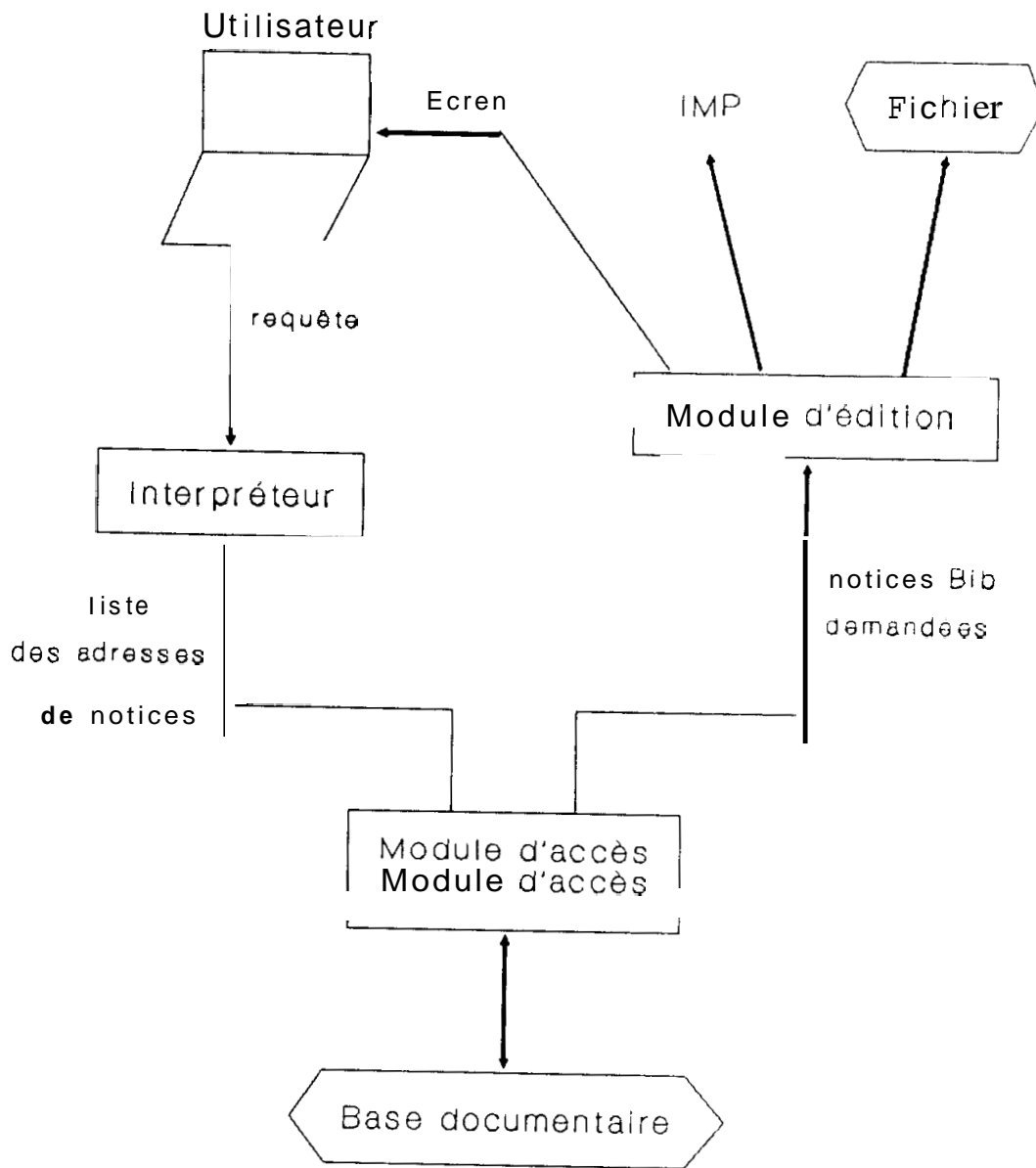
Une commande reconnue valide est exécutée par l'interpréteur. A ce niveau, le cheminement effectué par l'exécution d'une commande de recherche se schématise comme suit :



N<sub>i</sub> : nombre de references du mois i (i=1 a 12)  
 @i : adresse de rangement des adresses dans le fichier  
**INVERSE**

Schema d'exécution d'une commande de recherche

Le processus d'interrogation se schematise comme suit :



**Module d'interrogation de la base**

Annexe B  
Liste des mots vides utilisés par le module  
d'indexation du sous-système d'accès aux données  
volumineuses

A  
ABOUT  
ABOVE  
ACCORDING  
AFTER  
AGAINST  
ALMOST  
ALONE  
ALONG  
ALREADY  
ALSO  
ALTHOUGH  
ALWAYS  
AMONG  
AND  
ANOTHER  
ANY  
ANYONE  
ANYPLACE  
ANYTHING  
ANYTIME  
ANYWAY  
ANYWHERE  
ARE  
AROUND  
AS  
AT  
AWAY  
BACK  
BECAUSE  
BEFORE  
BELOW  
BETTER  
BETWEEN  
BOTH  
**BRIEFLY**  
BUT  
BY  
CLEARLY  
DOWN  
EACH

EARLY  
EASILY  
EITHER  
ENOUGH  
ESPECIALLY  
ETC  
EVEN  
EXCEPT  
FOR  
FROM  
FURTHER  
HARDLY  
HE  
HENCE  
HER  
HERE  
HERS  
HERSELF  
HIM  
HIMSELF  
HIS  
HOW  
HOWEVER  
I  
IF  
IN  
INTO  
IT  
ITS  
ITSELF  
LATE  
LEST  
MAINLY  
ME  
MERELY  
MINE  
MOST  
MUCH  
MY  
MYSELF  
NEARLY  
NECESSARY  
NEVER  
NEVERTHELESS  
NO  
OF  
OFF  
OFTEN  
ON

ONCE  
ONE  
ONES  
ONLY  
OR  
OTHER  
OTHERS  
OUR  
OURS  
OURSELVES  
OUT  
OUTSIDE  
OVER  
PER  
PROVIDED  
PROVIDING  
QUICKLY  
QUITE  
RATHER  
REGARDING  
RELATIVELY  
ROUND  
SELDOM  
SHE  
SINCE  
SLOWLY  
SOME  
SOMETIMES  
SOON  
STILL  
THAN  
THAT  
THE  
THEIR  
THEIRS  
THEM  
THEMSELVES  
THERE  
THEREBY  
THEREFORE  
THEY  
THOUGH  
THROUGH  
THROUGHOUT  
THUS  
TILL  
TO  
TODAY  
TOGETHER  
TOMORROW

TOO  
TOWARD  
TOWARDS  
UNLESS  
UNTIL  
UP  
UPSTAIRS  
US  
USUALLY  
WE  
WHATEVER  
WHEN  
WHENCE  
WHENEVER  
WHERE  
WHEREAS  
WHETHER  
WHICH  
WHILE  
WHO  
WHY  
WILL  
WITH  
WITHIN  
WITHOUT  
YESTERDAY  
YET  
YOU  
YOUR  
YOURS  
YOURSELF  
YOURSELVES