

# PROBLEMATIQUE GENERALE DE LA RECHERCHE D'INFORMATION

## Notions de base et typologies de la recherche d'information

### Questions de terminologie

### Typologies des recherches d'information

## Indexation humaine / indexation automatisée

### Notions sur l'indexation automatisée

## Principes universels de la recherche d'information

### Notions de base de la RDI

### Les opérateurs de recherche

## Différences entre la recherche sur Internet et la RDI

---

### Notions de base et typologies de la recherche d'information

Pour la recherche d'information, nécessaire de bien connaître :

- les types d'informations : **primaire et secondaire**
- la nature des informations recherchées
- les modes de recherche
- le contexte de la recherche
- les principes de base de la recherche d'information
- les sources
- les outils de recherche

---

### Questions de terminologie

- **Recherche documentaire :**

Définition officielle :

- **Action, méthodes et procédures ayant pour objet de retrouver dans des fonds documentaires les références des documents pertinents (Vocabulaire de la documentation, AFNOR, 1987). Ensemble des techniques et modalités permettant de sélectionner l'information dans un fonds documentaire structuré en fonction de critères de recherches propres à l'utilisateur.**
- **Recherche Documentaire Informatisée :**  
(ou RD Automatisée)  
Recherche documentaire utilisant un logiciel documentaire, sur ordinateur ; implique l'élaboration d'une démarche, d'une stratégie de recherche :
  - définition des mots-clés, des clés d'accès...
  - élaboration d'une équation de recherche...>> RDI plutôt assimilée à l'interrogation des banques de données.
- **Recherche de l'information :**
- Définition officielle :
  - **Action, méthodes et procédures ayant pour objet d'extraire d'un ensemble de documents les informations voulues (d'après l'AFNOR, 1979). Dans un sens plus large, toute opération (ou ensemble d'opérations) ayant pour objet la recherche, la collecte et l'exploitation d'informations en réponse à une question sur un sujet précis.**

### **Quels problèmes aujourd'hui pour ces définitions ?**

Toutes ces définitions sont bousculées par Internet, qui permet à la fois :

- de rechercher des références de documents : recherche bibliographique ou recherche documentaire classique ("secondaire")
- de rechercher des documents entiers, sous forme numérique : recherche de type "primaire"
- de rechercher des informations (exploitation des documents)

>> ce qui explique que l'on parle surtout de **recherche de l'information sur Internet**, pour signifier la numérisation du document et la primauté accordée désormais à l'information sur le document, au contenu sur le contenant



## Typologies des recherches d'informations

Il existe une grande diversité de modalités, de stratégies et de types de recherches d'informations : recherches en accès direct, recherches hypertextuelles, recherches par requêtes, recherches multicritères, recherche factuelle, etc. On peut distinguer les recherches d'information notamment selon :

- **la nature des informations** recherchées : qu'est-ce qui est recherché ?
- **les modalités de la recherche**, les modes d'accès : comment rechercher ?

### A/ Typologie selon la nature des informations recherchées :

- **Recherche factuelle :**  
Recherche d'informations précises, en réponse à une question précise. Recherche d'informations structurées, numériques ou alphanumériques, généralement contenues dans des banques de données, des annuaires, des répertoires ; ou bien recherche d'informations non structurées, dans une base textuelle.
- **Recherche documentaire :**  
recherche de documents à partir de leurs références ; recherche d'informations secondaires, aboutissant au document primaire.
- **Recherche bibliographique :**  
sous-ensemble de la recherche documentaire, limitée aux références bibliographiques des documents ou aux adresses des sites web.
- **Recherche contextuelle sur le contenu :**  
recherche à partir d'un mot ou d'un groupe de mots, pour

aboutir à un texte (à l'information primaire). Recherche en texte intégral.

## B/ Typologie selon les différents modes de recherche :

Cette typologie est plus importante que la précédente, car les modes d'accès induisent des compétences, des comportements, des pratiques et des résultats très différenciés. Par ailleurs, la recherche d'information sur Internet permet aujourd'hui de combiner tous les modes de recherche, auparavant très distincts, voire exclusifs les uns des autres.

Le tableau ci-dessous cherche à présenter une vue d'ensemble de ces différents modes de recherche, résumés ici à quatre grandes modalités, correspondant chacune à des principes de recherche différents, à des figures ou des métaphores du savoir distinctes, à des "gestes cognitifs" ou des démarches intellectuelles également distinctes.

### LES QUATRE PRINCIPAUX MODES DE RECHERCHES D'INFORMATIONS

<i>Modes de recherche</i>	<i>Principe, démarche intellectuelles</i>	<i>Type d'information concernée en priorité</i>	<i>Exemples d'outils</i>
<b>Recherche par navigation arborescente</b>	<b>Figure de l'Arbre</b> Démarche <b>systematique</b> , du général au particulier  Recherche <b>par menus successifs</b>	<b>Information structurée</b> , organisée en plan de classement  Information <b>secondaire</b>	<b>Tables des matières des livres</b>  <b>Classifications documentaires</b> (CDU, Dewey)  <b>Annuaire thématiques du web</b> (Yahoo, Nomade...)  <b>Page d'accueil d'un site web</b>
<b>Recherche par</b>	<b>Figure du Réseau</b>	<b>Information non</b>	<b>Renvois dans une</b>

<b>navigation hypertextuelle</b>	Démarche <b>associative</b> , d'une notion à l'autre.  Navigation dans un réseau de noeuds et de liens	<b>structurée</b>  Information <b>primaire</b> , brute	encyclopédie  <b>Hypertextes</b> sur CD-ROM  <b>Sites web</b>
<b>Recherche par requête sur les métadonnées du document</b>	<b>Principe de l'Index</b>  Démarche <b>d'indexation</b> de l'information  Recherche <b>par champs, logique booléenne</b>	Information <b>structurée en champs</b> .  Information <b>secondaire</b>	<b>Index des livres</b>  <b>Banques de données</b>  <b>Catalogues de bibliothèques</b>
<b>Recherche par requête sur le texte intégral</b>	<b>Texte</b>  Démarche <b>d'analyse linguistique</b>  Recherche <b>contextuelle sur le contenu</b>	<b>Information non structurée</b>  Information brute, <b>primaire</b>	Outils de <b>TALN</b> ( <i>Traitement Automatique du Langage Naturel</i> )  <b>Moteurs de recherche</b>  <b>Outils linguistiques</b>

Pour simplifier, **deux grandes méthodes** de recherche documentaire, actuellement juxtaposées :

- **la recherche dans les bases de données, c.a.d. dans des champs structurés :**
  - > recherche booléenne, par mots-clés
  - la "RDI" (Recherche Documentaire Informatisée) classique**
- **la recherche dans les corpus de textes, dans le texte intégral :**
  - > recherche linguistique
  - la recherche d'information sur Internet**

Mutations accélérées des méthodes de recherche et d'accès avec Internet : essor des recherches linguistiques.



---

## Indexation humaine / indexation automatisée

Pour l'analyse et la description des documents, deux grandes réponses apportées par la documentation :

- **l'indexation par un langage documentaire**
- **l'indexation automatisée du texte intégral**

- **Langage documentaire :**

langage artificiel élaboré pour éliminer les problèmes d'ambiguïté du langage naturel. Un langage documentaire vise à fournir une représentation univoque et formalisée des documents.

Exemples de langages documentaires : classifications décimales (Dewey, CDU...), listes de vedettes-matières (RAMEAU), thésaurus.

- **Indexation automatisée :**

traitement linguistique du texte intégral.

## Notions sur l'indexation automatisée

L'indexation automatisée repose sur les techniques de **TALN** : Traitement Automatique du Langage Naturel.

Au fondement de la recherche d'information sur Internet, notamment dans les moteurs de recherche.

Plusieurs niveaux d'analyse du texte intégral :

- **niveau morphologique** : reconnaissance du mot
- **niveau lexical** : réduction du mot à sa forme canonique > lemmatisation
- **niveau syntaxique** : niveau d'utilisation de la grammaire
- **niveau sémantique** : niveau de la reconnaissance des concepts

**Sur Internet : l'indexation morphologique est généralement le seul niveau d'indexation utilisé par les moteurs de recherche. La**

**plupart des moteurs n'éliminent pas les mots-vides et prennent en compte tous les mots.**

- **Indexation en texte intégral ou plein texte (full text) :**

- Fondée sur l'analyse morphologique des mots : leur **forme**.

- Méthode fondée sur **l'élimination des mots-vides** (à partir d'un dictionnaire de termes) et la **constitution d'un index des termes non éliminés**, considérés comme des chaînes de caractères.

- Recherche se fait selon logique booléenne ou avec des opérateurs de proximité ; méthode de recherche en full text quotidiennement utilisée dans les banques de données, pour les recherches portant sur les chaînes de caractères, dans les résumés par exemple.

Exemple : dans la phrase "*Prolétaires de tous les pays : unissez-vous*", l'indexation en full text éliminera les termes : de, tous, les, vous, et gardera "prolétaires", "pays" et "unissez"

A la recherche, il suffira de taper l'un de ces termes, ou une combinaison des termes, pour retrouver la phrase.

### **Quels problèmes et quels défauts de l'indexation morphologique ?**

Inconvénients évidents de cette méthode très fruste d'indexation :

- tous les mots non vides mis sur le même plan :

- pas de prise en compte de l'ordre des mots

- apparition des différentes formes d'un mot : par ex. un verbe va apparaître plusieurs fois sous des formes différentes

- l'analyse porte seulement sur des mots isolés (des unitermes), et délaisse toutes les expressions (les syntagmes), souvent porteurs de sens :

- ex. : pomme de terre donnera deux mots "pomme" et "terre", analysés séparément

- polysémie, synonymie du langage naturel pas prise en compte :

- > vol = aussi bien vol d'avion que vol à la tire

- stricte équivalence entre termes de recherche et termes indexés : pas de faute de frappe...

>> limites très sérieuses de l'analyse morphologique, qui peut générer beaucoup de « bruit ou de silence documentaire »



---

## Principes de la recherche d'information

Deux notions transversales de la recherche d'information (surtout dans la "RDI") :

**le bruit et le silence documentaire :**

- **bruit** : documents retrouvés non pertinents
  - **silence** : documents pertinents non retrouvés
- > ces deux indices définissent la **pertinence** (*relevance*) d'un système documentaire :
- 

## Notions de base de la RDI

- **Structuration de l'information :**

- une notice est découpée en **champs**
- les champs sont **interrogeables** (indexés) ou **non interrogeables** (non indexés)
- les champs indexés sont interrogeables par des **mots-clés**
- pour l'interrogation d'une base de données, utilisation d'opérateurs de recherche

- **Notion de zones, champs ou rubriques :**

Toutes les informations présentes dans un document sont regroupées dans des **champs** :

- appelés aussi zones ou rubriques
- chaque élément de description sera isolé dans un espace précis, identifiable : un champ.

Ex. : champ Auteur, Titre, Editeur...  
Ensemble des champs définis constitue la « structure » de la base de données, ou du moins de l'enregistrement, de la notice.

- **Notice, enregistrement :**

**Notice ou enregistrement, ou fiche, ou référence = champs + données**

- **Les clés de recherche : les mots-clés**

Une fois saisies les informations, le problème posé est celui de leur recherche, de leur accès.

Deux possibilités :

- **accès séquentiel** (par balayage : l'ordinateur devrait lire toutes les notices une par une...) : > lenteur inacceptable !

- **accès direct** : spécificité des systèmes de recherche documentaire, possibilité d'accéder directement à une notice, par l'intermédiaire d'une clé de recherche, contenue dans un index

> notion de « **mot-clé** » (*keyword*)

Tous les mots-clés sont regroupés dans des fichiers séparés des notices : les **fichiers index** (parfois appelés lexiques selon les logiciels)

Ex. : index des auteurs, des titres, des dates de publication...

---

**Les opérateurs de recherche (booléens, troncature, proximité...)**

**TABLEAU DES OPERATEURS DE RECHERCHE**

<i>Catégories d'opérateurs</i>	<i>Définition</i>	<i>Symboles</i>	<i>Types</i>	<i>Fonction</i>	<i>Exemples</i>
--------------------------------	-------------------	-----------------	--------------	-----------------	-----------------

<b>Troncature</b> <b>Wildcards,</b> <b>Wildcard</b>	Substitution d'un symbole à des caractères	<ul style="list-style-type: none"> <li>• *</li> <li>• +</li> <li>• ?</li> </ul>	<ul style="list-style-type: none"> <li>• <b>Troncature à droite</b></li> </ul>	Recherche sur tous les mots contenant la même racine ou le même préfixe	franco* : >> francophone, francophonie, francophobe
			<ul style="list-style-type: none"> <li>• <b>Troncature à gauche</b></li> </ul>	Recherche à partir d'un suffixe	*phobe : >> technophobe, agoraphobe, etc.
		<ul style="list-style-type: none"> <li>• ?</li> <li>• #</li> </ul>	<ul style="list-style-type: none"> <li>• <b>Masque ou troncature centrale</b></li> </ul>	Remplace un ou plusieurs caractères dans un mot.	francopho?e > francophobe et francophone
<b>Opérateurs booléens</b>	Logique booléenne, permettant des relations entre les termes de la recherche	<ul style="list-style-type: none"> <li>• <b>ET</b></li> <li>• <b>AND</b></li> <li>• +</li> </ul>	<ul style="list-style-type: none"> <li>• <b>Opérateur d'intersection</b></li> </ul>	Relie deux mots-clés se trouvant ensemble dans le même document	Pêche ET Bretagne :  >> les documents traitant de la pêche en Bretagne
		<ul style="list-style-type: none"> <li>• <b>OU</b></li> <li>• <b>OR</b></li> </ul>	<ul style="list-style-type: none"> <li>• <b>Opérateur d'union</b></li> </ul>	Relie deux mots-clés dont l'un ou l'autre, ou les deux doivent se trouver dans le document	Pêche OU Bretagne :  >> les documents traitant ou bien de pêche ou bien de Bretagne, ou des deux à la fois
		<ul style="list-style-type: none"> <li>• <b>SAUF</b></li> <li>• <b>NOT</b></li> <li>• -</li> </ul>	<ul style="list-style-type: none"> <li>• <b>Opérateur d'exclusion</b></li> </ul>	Relie deux mots-clés dont le premier doit être présent et le second absent dans un document	Pêche SAUF Bretagne :  >> les documents traitant de la pêche sauf en Bretagne
<b>Opérateurs numériques</b>	Permettent des recherches selon <b>des critères quantitatifs</b>	<ul style="list-style-type: none"> <li>- égal : =</li> <li>- supérieur : &gt;</li> </ul>		Recherche sur les dates, par exemple	<ul style="list-style-type: none"> <li>• = 1995 : en 1995</li> <li>• &gt; 1995 : depuis</li> </ul>

		<ul style="list-style-type: none"> <li>- supérieur ou égal : &gt;=</li> <li>- inférieur : &lt;</li> <li>- inférieur ou égal : &lt;=</li> <li>- intervalle entre deux dates : :</li> </ul>			<p>1995</p> <ul style="list-style-type: none"> <li>• &gt;= 1995 : depuis ou en 1995</li> <li>• &lt;= 1995 : avant ou en 1995</li> <li>• 1991:1995 : de 1991 à 1995</li> </ul>
<b>opérateurs de proximité</b>	Permettent des recherches sur le texte intégral selon la proximité des termes	<ul style="list-style-type: none"> <li>• <b>ADJ</b></li> </ul>	<ul style="list-style-type: none"> <li>• <b>Opérateur d'adjacence</b></li> </ul>	Recherche sur des termes adjacents, dans l'ordre donné	<p>Fibre ADJ optique :</p> <p>&gt;&gt; texte contenant l'expression " fibre optique "</p>
		<ul style="list-style-type: none"> <li>• <b>nAV</b></li> </ul>	<ul style="list-style-type: none"> <li>• <b>Opérateur de distance</b></li> </ul>	Recherche sur des termes séparés par une distance <i>n</i>	<p>Ecole 1AV privée :</p> <p>&gt;&gt; école primaire privée, ou école technique privée</p>
		<ul style="list-style-type: none"> <li>• <b>NEAR</b></li> <li>• <b>M (pour Mot) ou W (Word)</b></li> </ul>	<ul style="list-style-type: none"> <li>• <b>Opérateur de présence</b></li> </ul>	Recherche sur des termes présents dans le texte, quelle que soit leur distance	<p>Fibre NEAR optique</p> <p>&gt;&gt; texte contenant les deux termes, mêmes séparés</p>

A noter : **l'universalité des opérateurs de recherche**, présents (partiellement ou en totalité) dans la totalité des systèmes de recherche d'information : catalogues de bibliothèques, bases et banques de données,

moteurs de recherche, outils d'indexation du texte intégral, annuaires thématiques...



### Différences entre la recherche sur Internet et la RDI "classique"

<b>Critères</b>	<b><i>Systèmes documentaires classiques</i></b>	<b><i>Internet</i></b>
<b>Caractéristiques générales de l'information</b>	<b>Homogénéité : des sources, des informations, des domaines, de la description, de la présentation, de la recherche...</b>	<b>Hétérogénéité ...</b>
<b>Nature de l'information</b>	<b>Informations secondaires</b>	<b>Informations primaires mélangées aux informations secondaires</b>
<b>Structuration de l'information</b>	<b>Information structurée, découpée en champs</b>	<b>Information plutôt non structurée. Utilisation réduite des métadonnées</b>
<b>Validité de l'information</b>	<b>Informations doublement validées : sources connues, responsabilité du professionnel de la documentation</b>	<b>Informations non validées : sources non connues, qualification sous la responsabilité de l'utilisateur</b>
<b>Cohérence de l'indexation</b>	<b>Indexation par un langage documentaire homogène, adapté à la base et au domaine</b>	<b>Pas de langage d'indexation.</b>
<b>Etendue de l'indexation</b>	<b>Indexation de la totalité de la base</b>	<b>Indexation limitée du web visible (30 %), web "invisible" non indexé...</b>

Niveau de l'indexation	Indexation humaine (niveau sémantique) ou automatisée (plusieurs niveaux)	Indexation automatisée, limitée souvent au plus bas niveau (analyse morphologique) ; mais développement de l'indexation lexicale et syntaxique.
Représentation des documents	<b>Normalisée</b> : notices textuelles, de description bibliographique. Format unique de représentation des informations	<b>Non normalisée</b> : multiplicité des formats, liberté des auteurs... Exemple des 3 formes possibles d'un article : - une seule page HTML - plusieurs pages HTML - " chapeau " + document en format PDF (ou DOC)
Information accessible en recherche	<b>Notices</b> structurées, portant sur les différentes caractéristiques des documents (auteurs, titre, mots-clés...), et normalisées selon le système documentaire	Dans les annuaires : <b>Adresses de sites.</b>  Dans les moteurs : <b>texte des pages HTML, critères de date</b>
Volume d'information	<b>Limité</b> (qqs centaines de milliers de documents), maîtrisé par le responsable du système. Croissance linéaire de la base	<b>Très grande quantité d'informations ; volume redondant, non maîtrisé ;</b> pas de filtrage. Croissance exponentielle du web
Enrichissement de l'information	Ajout de documents, mais <b>stabilité du contenu et des notices des documents</b>	Information <b>très changeante</b> , modifications constantes du contenu et de l'organisation des sites
Responsabilité de l'alimentation du système	<b>Base documentaire alimentée par les professionnels de l'information : souci de cohérence de l'information</b>	<b>Sites web alimentés par les webmasters : souci de présentation, de diffusion et d'impact de l'information</b>
Outils de recherche	<b>Outil d'interrogation unique (logiciel documentaire),</b>	<b>Diversité des outils de recherche, à adapter selon</b>

	<b>adapté à la base et au langage documentaire.</b>	<b>les types de recherche. Mélange de tous les modes de recherche.</b>
<b>Compétences ou connaissances requises pour l'utilisateur</b>	Connaissance minimale : - <b>du domaine documentaire</b> - <b>de la structure des notices, des mots-clés</b> (et du langage d'indexation utilisé)	Connaissance : - <b>des sites pertinents</b> > l'offre - <b>des outils de recherche</b> - <b>des méthodes de recherche</b>
<b>Qualité de la recherche</b>	Pour une recherche par mots-clés :  - dépend directement de la <b>qualité de l'indexation manuelle</b>  - <b>aucune intelligence de l'outil informatique</b> : retrouve <b>exactement</b> les notices contenant les mots-clés demandés	Qualité de recherche : <b>aléatoire.</b> Dépend : - <b>de la couverture du domaine</b> par des sites accessibles en recherche (> inégalité des domaines couverts sur Internet) - <b>de l'expertise de l'utilisateur</b> - <b>des performances de l'outil</b>  <b>Intelligence distribuée dans les outils</b> ; recherche retourne la question qui <b>ressemble</b> le plus à la question posée > recherche par approximations successives... Notion de " <b>serendipité</b> " ou " <i>découverte, par chance ou sagacité, de résultats que l'on ne cherchait pas</i> ".
<b>Multilinguisme</b>	<b>Pas à la charge de l'utilisateur</b> : mots-clés dans les bdd en une ou deux langues, même si documents en plusieurs langues. > pas de prise en compte de la langue des documents dans la recherche	<b>A la charge de l'utilisateur</b> : - choix de la langue des documents dans le moteur, tri...



---

**Date de dernière mise à jour :  
31 janvier 2004**

***Ce support de cours peut être librement exploité, sous réserve de  
citer son origine.***

**© URFIST Bretagne-Pays de Loire, Alexandre Serres, 2002**

---