

Information Retrieval - Introduction

Recherche et Extraction
d'Information

Anne-Laure Ligozat/Xavier Tannier



Index

[illegible]

J'ai de la chance

A

Abtchoud S. 799, 945
Abtchri J. R. 421
accès aux données distantes
accès direct (index) 345
accès séquentiel (accès) 446
accès séquentiel (physique) — voir séquençage physique
accepter une préface 911
Achilles S. 690
Adiba M. E. 575, 734
Adelman L. 500, 510
administrateur de la base de données — voir DBA
administrateur des données 15
autoajuste disposé — voir disposition
afféction relationnelle 158
réindus cible 569
agent 707
Agrawal R. 661, 841, 945, 948
Aha A. V. 382, 392, 824, 946
algèbre relationnelle 149
implantation 603
objet 180
opérations primitives 190, 203
règles de transformation 181, 592
algorithme de chasse 393
algorithme de duplication de Codd 236
ALG (SQL) — voir dialecte
Allen P. W. 424
ALPHA — voir DSL ALPHA
ALTER DOMAIN (SQL) 258
ALTER TABLE (SQL) 100, 261
Alman E. B. 873
American National Standards Institute — voir ANSI
analyse informatique 456, 465
Anderson E. 665
annonces de nalsc à jour 345, 329, 356, 378
ANSI 73
ANSI/SPARC 33, 57
ANSI/X3 57
ANSI/X3/SPARC Study Group on Data Base Management Systems — voir ANSI/SPARC
Anton J. 664
appel des procédures distantes — voir RPC
APPEND (QUEL) 496
équivalences en lignes 9
arbre de requête 358
arbre de syntaxe abstraite — voir arbre de requête
arbre de recherche analytique — voir tree
architecture ANSI/SPARC 33
en SQL 80
ARIES 463
arriver au jour digre
Armstrong W. W. 326, 328
Arps M. 661
Arsenault R. L. 92
assertion (SQL) — voir CREATE ASSERTION
association 12, 405, 406, 410
CQ 390
récurative 414
association (RM/T) 426
associativité 170
Astroman M. M. 295, 296, 873
Atkinson M. P. 802, 822
atomisme
relations 353
transactions 426, 441
valeurs xéniques 63, 164, 642
attribut 98, 97
authentification — voir mot de passe
automatique inné 700
automatique — voir sécurité
auxiliaire 427
AVG — voir fonction d'aggrégation
axiome 906
axiome d'échafaut 525
axiome de base 600, 923
axiomes de Armstrong 320, 328

B

B-trees 860
Badal D. Z. 544
Badriloff F. 832, 944, 945, 947
Baecker J. 800
Barney C. M. 872
Barney M. F. 884
base de connaissances 302
base de données 3, 10
statistiques 16
702, 890
base de données déductive 910
base de données distribuée 55, 606
pétence fonctionnelle 439
base de données experte 940
base de données extentionnelle 923
base de données interrelationnelle 925
base de données logique 840
base de données relationnelle 110
base de données statistique 508
Bates J. C. 43, 45, 874
Bayer R. 478, 851, 875
BCJR 237, 354, 361
BLA 10
Bledley D. A. 884
Blewett D. 825
Blewett C. 392, 392, 393, 946
BLOOM DECLARE SECTION (SQL) 283
BEGIN TRANSACTION 439
Bell D. 729
Bentley J. L. 454
Berenson F. A. 683, 419, 473, 534, 576, 729
Bird (DDI) 884, 895
Bittan D. 437
Björner A. 800
Björner A. 800

Information retrieval actors

Collection :

a set of
documents



User :

an information
need and/or
a task to perform



IR system: the tool that should
find relevant documents for
the user information need



Information Retrieval



Web Images Videos Wikipedia More ▶

librairie orsay

Search

Advanced Search

Home > Web results 1-10 of 119,601 for language:fr librairie orsay, Page 1 - Next page



Librairie Orsay, Livre - Recherche adresse

Bonnes adresses > Essonne > **Orsay** > Culture, Loisirs et Voyages > **Librairie Librairie Orsay** Résultats 1 à 1 sur 1 pour **Librairie Orsay** Librairie Du Lycee Donnez votre avis 0% 57 Rue Paris - 91400

www.justacote.com/orsay-91400/librairie

Cached - Bookmark



Librairie Orsay - librairies

Culture high tech **librairie Orsay** - Toutes les infos sur Culture high tech **librairie Orsay** - Avis des internautes, téléphone, horaires, itinéraires, adresse et plan

fr.nomao.com/orsay/faire-du-shopping/librairie.html

Cached - Bookmark



Achat - Vente Librairie Orsay - 91400, Cession Librairie Orsay - 91400

Des milliers d'annonces de **Librairie Orsay** - 91400 à vendre ou à céder avec Vivastreet **Orsay** - 91400, trouvez votre **Librairie** parmi plus de 11 000 ANNONCES 100% GRATUITES Achat - Vente Cession

fonds-commerce.vivastreet.fr/annonces-commerce-pas-de-porte-orsay-91400/q/librairie

Cached - Bookmark

Site type:

- » [Blog](#)
- » [Forum](#)

Multimedia:

- » [Video](#)

Filetype:

- » [pdf](#)
- » [swf](#)
- » [text](#)
- » [word](#)

Related terms:

- » [Art Contemporain](#)
- » [Art Moderne](#)
- » [Assemblée nationale](#)
- » [Centre Pompidou](#)
- » [Grand Palais](#)
- » [Musée d'Orsay](#)
- » [Musée National](#)

Information Retrieval

- Where is the bookstore closest to home?
- Who is presently leader of the rugby Top 14?
- Quels sont les titres mentionnés à la une du journal Le Monde d'aujourd'hui ?
- Que rapporte la une du Monde d'aujourd'hui sur la politique étrangère ?
- Quels sont les films qui passent ce soir sur la TNT ?
- Dans quels films Jean Rochefort et Philippe Noiret ont-ils joué ensemble ?
- Quels sont les logiciels d'installation de logiciels sous Linux/Debian ?
- Comment peut-on installer des logiciels sous Linux/Debian ?
- What is the English word for "givre" ?
- Who was Claude Bernard?

Questions

- What kind of results are expected?
- How to evaluate the results relevance?
- How to formulate a query?
- ...

Information vs. données

"Les **données** sont reçues, stockées et retrouvées par un endosystème. Les données sont impersonnelles ; elles sont disponibles pour tout utilisateur du système.

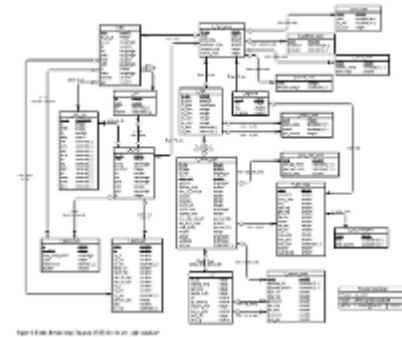
L'**information**, en revanche, est un ensemble de données qui correspond à un besoin particulier.

Le concept d'information a des composantes personnelles et temporelles absentes du concept de donnée."

(R. R. Korfhage, 1997)

Diversity of information needs (1/2)

- Search for a **known element**
- The user knows exactly which elements they look for
 - and can recognize them
 - *Example:* search for a bibliography item
 - **Databases (SQL, XQuery, etc.)**
- Search for a **general information**
 - The user searches for information about a subject
 - and may not recognize relevant information
 - or information can be partially relevant
 - *Example:* reforms of research in France
 - **Traditional information retrieval**



Diversity of information needs (2/2)

- Search for **precise information**

The user searched for specific information
but ignores under which form it is

- Partial answer not relevant
- *Example* : What day was president *Kennedy* killed?

➤ **Information extraction and question answering systems**



- **Exploration**

The goal is not to answer a particular question,
but to go through data to discover information
about a subject or a domain

➤ **Navigation**



Diversity of information sources



- **Location** of information
 - Local or distant resources
 - Problems : disponibility, identification, distribution on several sources, format variability (character encoding and content description)
- **Nature** of the resource files
 - **Databases**: well described formats, non ambiguous query languages (ex : SQL for relational databases)
 - **Annotated files**: more or less described formats, presentation and/or semantic description annotations, query languages (ex : XSLT/XPath for XML files)
 - **Text files**: few or not described formats, known language(s) or not, more or less regularity between documents of a same class, no generic interpretation possible (NLP probleme)

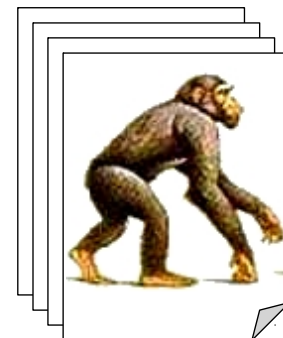
Diversity of problems

- Difficulty to **access, cover and treat**
- Document bases **very large**, distributed on
 - **numerous supports** in **different places**
- Difficulty to define **relevance**
 - How does a document respond to the **user information need**?
 - What is the **relevance**? How to measure it?
- Difficulty to **exploit**
 - Relevant documents may not be in the query language
 - The desired information may not be easy to identify in the document.



IR main evolutions

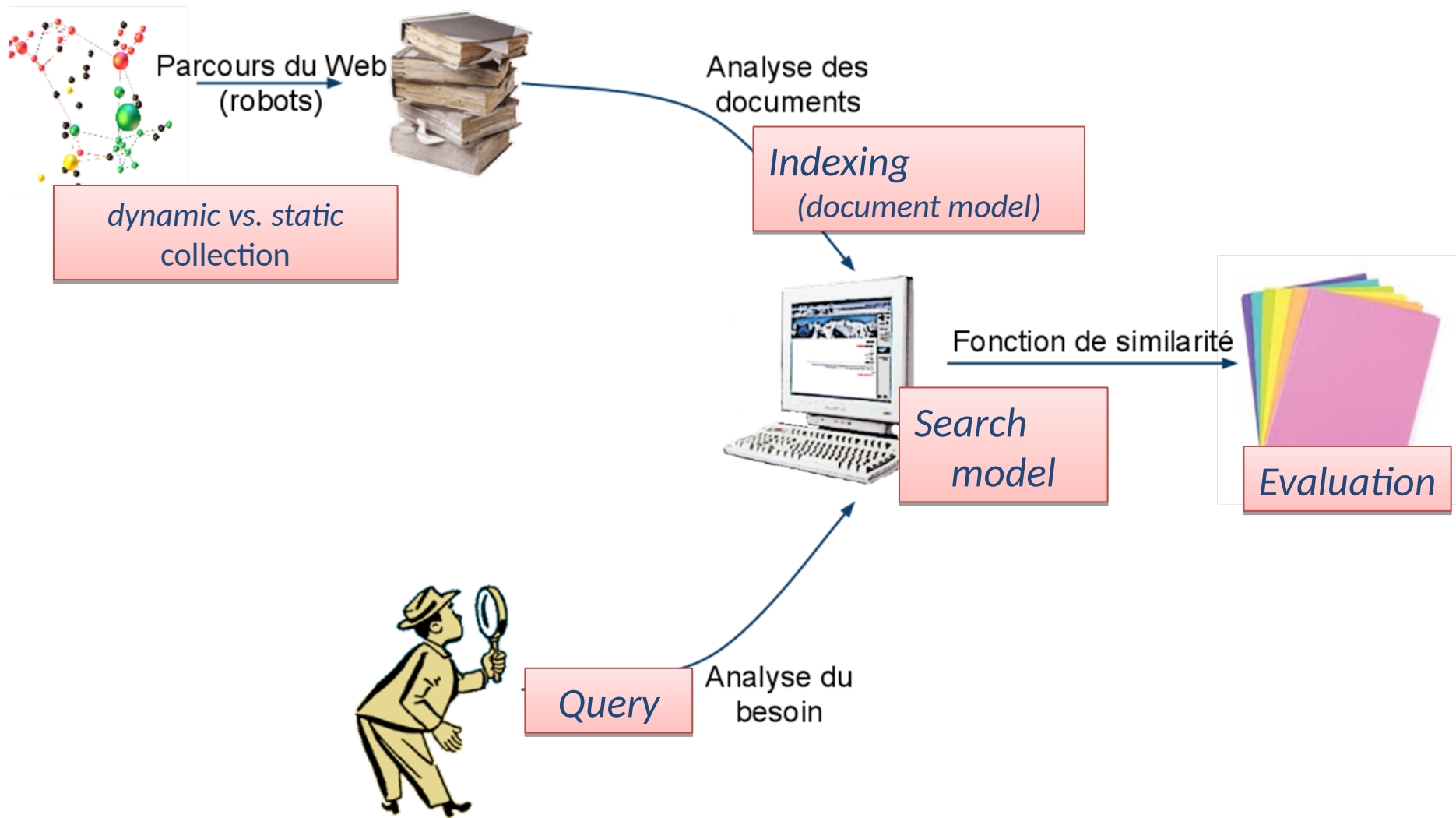
- In the beginning
 - Small and structured document bases
 - Acces by metadata, and rarely by whole text
 - Use of specific languages by specialists
- Nowadays
 - Digital multimedia documents
 - Various representation formats (raw text, HTML, XML, PDF, RTF,...)
 - More and more unstructured data
 - Huge quantity (Web...)
 - Social component



Web IR

- On the Web: massive use by non expert users
 - domain with major economic importance
 - typical query: a few keywords
 - users adapt to existing tools
- Part of the Web is not directly accessible (invisible web, including restricted access pages and dynamic pages)
- Strongly multilingual information: the documents containing query information may be in different languages
- Information not always reliable
- Information visualization particularly important: result ranking, extract presentation, relevant segment extraction etc.

Information Retrieval



IR difficulties

Humans and languages

Les difficultés de la RI : le facteur humain

- Le besoin d'information de l'utilisateur est parfois **vague** et toujours **subjectif**.
 - La **perte d'information** entre la réalité du besoin d'information et son expression peut être importante.
 - La pertinence d'un document pour une requête est une notion variable et très complexe à définir.
- ⇒ Il ne peut pas exister de système de recherche d'information parfait
- L'évaluation d'un système dépasse les aspects habituels de performance informatique
- L'humain est subjectif, versatile, et il utilise un **langage "naturel"** !



Les difficultés de la RI : le facteur "langage"

- À la différence des langages artificiels, le **langage "naturel"** est :
 - **Implicite** : tout n'est pas dit dans les textes et leur compréhension requiert une importante connaissance sur le contexte et sur le monde
 - **Redondant** : la langue offre de nombreuses façons de formuler le même contenu
 - **Ambigu** : un même énoncé peut souvent être interprété de différentes façons
 - La recherche d'information est encore compliquée par le fait que :
 - Les mots peuvent jouer des rôles différents dans les textes
 - Les atomes de sens peuvent être des mots ou des groupes de mots (termes)
- Il est compliqué de **formuler son besoin d'information**
(perte d'information entre besoin et requête)



Caractère implicite de la langue

- Connaissance du langage et des **conventions langagières**

Q : *Le voisin est-il chez lui ?*

R : *Sa voiture est devant le portail*

(implicature conversationnelle)

- Connaissance du **contexte**

C'est la deuxième fois qu'il reçoit un carton

(Sport ? Courrier ? Accident ?)

- Connaissance du **monde**

La Nouvelle-Zélande va tailler la France en pièces.

(métonymie + langage figuré + actualité du rugby)

- **Déduction** (présupposition)

Ravaillac a assassiné Henri IV en 1610.

⇒ Henri IV est mort en 1610.

Caractère *redondant* de la langue

- Au niveau lexical
 - **Synonymie** : vélo et *bicyclette*
 - **Hyperonymie** et **hyponymie** : véhicule \triangleleft vélo \triangleleft VTT
 - **Méronymie** et **holonymie** : pédale \diamond pédalier \diamond vélo
- Abréviations et sigles
 - *S'il-vous-plaît* et *SVP*, VTT et Vélo Tout Terrain
- Entre mots et expressions
 - **Périphrases** : *lave-vaisselle* et *machine à laver la vaisselle*
 - **Définitions** : *selle* et *petit siège, le plus souvent de cuir, d'un cycle ou d'un véhicule à deux roues à moteur*
- Glissements de sens (synonymie contextuelle)
 - *Il a écrit un **papier/article** sur la recherche d'information*
 - **Vos **papiers/articles** s'il-vous-plaît !*



Caractère *redondant* de la langue

- La **paraphrase** (synonymie au niveau syntaxique)

- Qui sera élu par le peuple en 2022?
- Qui le peuple choisira-t-il dans 5 ans ?
- Qui sortira vainqueur des urnes pour le prochain quinquennat ?



- Synonymie et paraphrase ne sont pas transitives !

Paul ressort souvent **excité** de la **récréation**.
Paul profite bien de la **distraction** de la récréation.
Entre deux cours, les **loisirs** sont bénéfiques à Paul.
Une période de **détente** entre deux cours ne fait pas de mal à Paul.
Paul se **repose** entre deux cours.



Caractère *ambigu* de la langue

Les **homonymes** sont des mots qui ont une même graphie mais des sens différents



WIKIPÉDIA
L'encyclopédie libre

[Accueil](#)
[Portails thématiques](#)
[Index alphabétique](#)
[Un article au hasard](#)
[Contacter Wikipédia](#)

▼ [Contribuer](#)
[Aide](#)
[Communauté](#)
[Modifications récentes](#)
[Accueil des nouveaux arrivants](#)
[Faire un don](#)
► [Imprimer / exporter](#)
► [Boîte à outils](#)

Article [Discussion](#)

Lire [Modifier](#) [Afficher l'historique](#)

Noyau

 Cette page d'*homonymie* répertorie les différents sujets et articles partageant un même nom.

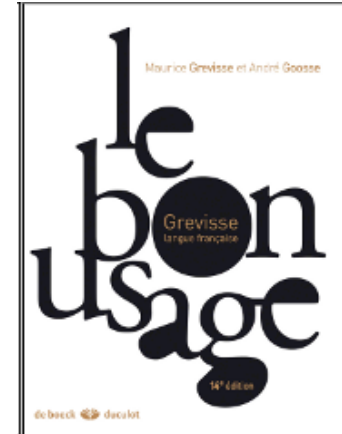
De manière générale, un **noyau** est la partie centrale située au milieu d'un autre objet. Plus particulièrement, le terme peut faire référence à :

- en **biologie**, un **noyau** est un organite qui contient la plupart du matériel génétique ;
- en **linguistique**, un **noyau** est partie fondamentale du **syntagme**, entourée de ses **satellites** ;
- en **botanique**, un **noyau** est la partie centrale, dure, d'une **drupe** ou fruit à noyau ;
- en **électrotechnique**, un **noyau** est la pièce magnétique sur laquelle un fil conducteur est enroulé afin de réaliser une bobine ;
- en **fonderie**, un **noyau** est la partie d'un moule permettant la réalisation des parties creuses d'une pièce ;
- en **géologie**, un **noyau** est la partie centrale approximativement sphérique de la **Terre** ou d'une **planète** ;
- en **informatique**, un **noyau** (aussi appelé **kernel**) est la partie fondamentale de certains **systèmes d'exploitation** ;
- en **mathématiques**,
 - en **algèbre**, le **noyau** d'un **morphisme de groupes** est un sous-groupe particulier du groupe de départ,
 - en **analyse fonctionnelle**, un noyau est une fonction permettant de définir un **opérateur intégral** ;
- en **physique**, un **noyau** est la région centrale constituée des **nucléons** d'un **atome** ;
- en **bande dessinée**, **Noyau** est le nom de l'illustrateur **Yves Nussbaum** ;

Caractère *ambigu* de la langue

- Les **ambiguïtés syntaxiques** :

- Jean vend une tarte **aux pommes**.
- Jean vend une tarte **aux clients**.
- Jean a rapporté un vase **de Chine**.



- Les **anaphores** :

Ségolène trahit Martine. Son ancien mari lui en voulut longtemps.

- Les **ellipses** :

- Quelle est la taille **de la tour Eiffel** ? Le poids ?
- Les Stéphanois portent des **écharpes** vertes et les
- Toulousains des rouges et noires.

Mots composés

- Les **mots composés** sont beaucoup moins polysémiques
- Les rechercher ensemble dans les textes est bénéfique (mais compliqué)
- Ils ont un sens qui n'est pas la composition des sens des atomes
 - *Homme-grenouille*
 - *Pomme de terre*
 - *Traitement de texte*



© M. Heinrich, J. Negra

Morphologie

La **morphologie** est l'étude de la construction des mots (leurs structures, variations, et similitudes)

- L'**analyse morphologique** permet de décomposer un mot et d'extraire principalement :
 - La **racine**
 - Le **lemme**
 - La **catégorie morphosyntaxique** (catégorie grammaticale)
 - Les **traits morphologiques**

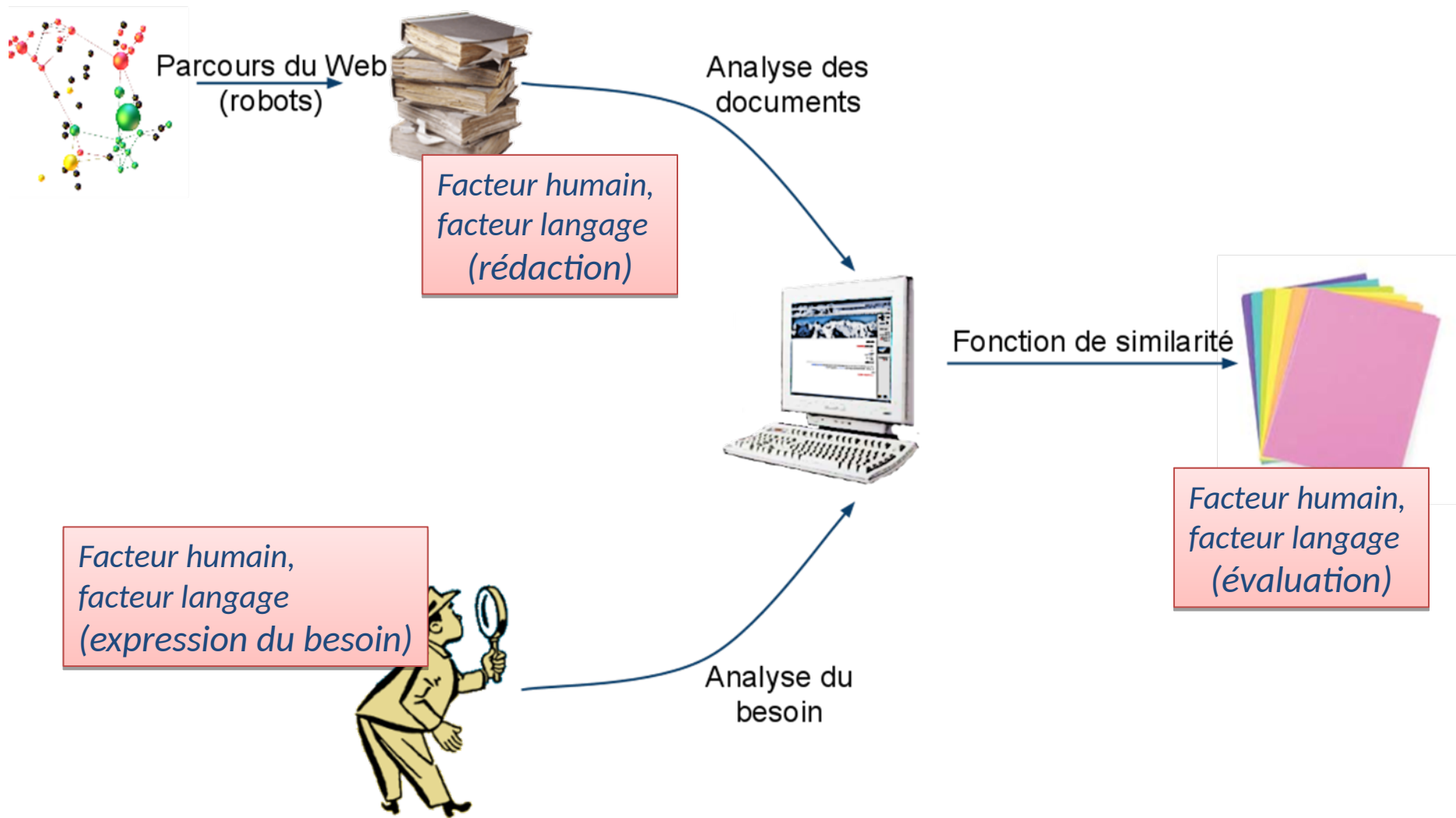
Construction de mots

Flexion

- **Composition**
- **Dérivation** (affixation)
 - **Préfixation**
 - **Suffixation**

(à combiner)

Recherche d'Information



Information Retrieval

Information retrieval is a statistical treatment of strings



- Enables to search large amount of information
- Applies to natural language text, does not require interventions from website creators, or particular representations of knowledge



- Machines do not understand information sense
- A search engine cannot perform inference of information cross-checking