

Quatrième Partie

4. L'accès à distance aux manuscrits

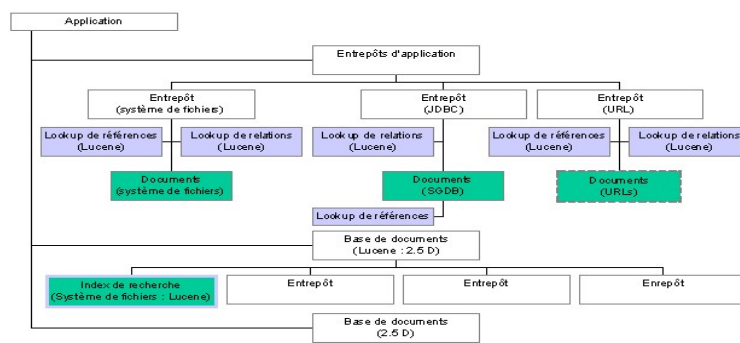


Figure n°. 60 : Aperçu global de SDX-2

4.1.1. La fonctionnalité de SDX

Il faut mentionner que SDX ne possède aucune interface graphique. C'est une application qui regroupe plusieurs autres applications développées en « open-sources ». Elle a son propre langage de programmation (XSP) pour unifier les commandements nécessaires à l'utilisation des différentes applications. Elle se compose de plusieurs logiciels qui permettent de *publier*, d'*indexer* et de *rechercher* des documents XML. Par exemple :

- ❑ **Cocoon**¹²⁰ : pour la publication, le SDX utilise le logiciel Cocoon. Cocoon est un outil qui aide à l'utilisation de XML et de XSLT pour des applications de serveur.
- ❑ **XML:DB** : SDX permet d'indexer des documents XML afin de faciliter leurs recherches. Pour indexer les documents, SDX utilise la base de données XML « XML:DB »
- ❑ **eXist** est une base de données qui propose un mécanisme d'indexation et de recherche en plein texte par mots clés.

«La force de SDX est donc de proposer une indexation complètement personnalisable qui, dans un cas extrême, pourrait se passer totalement du contenu du document original. On pourrait ainsi la réduire à des métadonnées par exemple ; c'est une pratique courante

¹²⁰ Pour savoir plus sur Cocoon voir l'adresse email : <http://xml.apache.org/cocoon/index.html>

chez le bibliothécaire et les archivistes qui n'indexent pas le contenu des documents dont ils ont la charge mais les moyens d'accéder à ces documents »¹²¹

4.1.2. *Le recherche des documents par SDX:*

Le logiciel SDX est structuré de manière à gérer que de l'information en XML. En conséquence, il peut indexer n'importe quel type de documents XML et il permet la configuration des champs de recherche de façon très souple bien adaptée aux informations gérées et aux besoins des utilisateurs. Par contre, toutes les bases de données qui n'existent pas en format XML peuvent être transformées rapidement en XML. La souplesse de SDX facilite une recherche simultanée dans plusieurs bases de documents, qu'elles soient définies auprès d'une même application ou dans différentes applications.

Pour la recherche des documents, il est possible sur SDX d'effectuer les deux types de recherche : une recherche textuelle ou par champ. « De plus, SDX offre non pas *une* mais *autant* d'indexations que l'on souhaite tout en se proposant d'offrir à l'utilisateur une syntaxe de recherche simple, à tout le moins plus simple que celle de XPath, SQL* ou autres langages de requête»¹²².

4.1.3. *L'architecture de SDX-1*

Pour la diffusion et la recherche d'une collection de documents XML, le SDX-1 comporte plusieurs composants¹²³ :

- Un système de stockage (un SGBD relationnel, tel que MySQL¹²⁴, InterBase¹²⁵, Oracle¹²⁶, Hypersonic SQL¹²⁷)

¹²¹ Brihaye, Pierrick. SDX-2 en tant qu'infrastructure générique pour la conception de systèmes documentaires : 19pages (P.7) (consulter 03/03/03) <http://sdx.culture.fr/sdx/journeeSDX/pres/brihaye/brihaye.html>

- ❑ Un serveur HTTP (comme celui de Tomcat, Apache ou autre)
- ❑ Un moteur de servlets Java (Tomcat ou autre)
- ❑ Un processeur XSLT pour transformer l'XML généré en HTML (Saxon)
- ❑ Un moteur de recherche (Lucene)
- ❑ Un générateur de pages XML (Cocoon-1)

Une application type se présentait en quatre volets :

- ❑ présentation des documents,
- ❑ navigation hiérarchique,
- ❑ recherche simple et avancée,
- ❑ Administration de la collection de documents.

Dans la version SDX-2 les ingénieurs du logiciel ont essayé de surmonter les inconvénients de la première version de SDX à plusieurs niveaux :

- ❑ Le servlet Cocoon pour améliorer la performance des XML, XSLT et XSP.
- ❑ Une autre amélioration est apporter à la fonction d'indexation, de recherche et d'affichage de document XML. Pour assurer cette fonction « il faut qu'un serveur Web, un moteur de servlet et SDX soient en marche et correctement configurés pour être prêts à servir les requêtes. La réalisation d'un tel environnement est appelée un *serveur* SDX, c'est-à-dire une instance du serveur en état de marche »¹²⁸
- ❑ Le moteur de recherche Lucene qui a été utilisé dans le SDX-1 est théoriquement remplaçable (à terme) mais pratiquement, peu de composants libres pourraient le remplacer. Un autre processeur XSLT peut être choisi sans risque

* SQL (Structured Query Language) voir http://www.w3schools.com/sql/sql_select.asp

¹²² Brihaye, Pierrick. SDX-2 en tant qu'infrastructure générique pour la conception de systèmes documentaires : Pistes de recherche et développement. <http://sdx.culture.fr/sdx/journeeSDX/pres/brihaye/brihaye.html#d0e28> dernière mise à jour le 06/02/2002

¹²³ Sévigny, Martin et Michel Botin. SDX en Quelques mots. <http://sdx.culture.fr/sdx/definition.html>

¹²⁴ <http://www.mysql.com/>

¹²⁵ <http://www.borland.com/interbase/>

¹²⁶ <http://www.oracle.com/>

¹²⁷ <http://hsqldb.sourceforge.net/>

¹²⁸ <http://www.nongnu.org/sdx/docs/html/doc-sdx2/fr/presentation/serveur.html>

d'incompatibilité, à la condition qu'il soit en Java (ex:Xalan). Xalan-Java est un processeur XSLT qui aide à transformer des documents XML en HTML, texte ou autres types de documents XML « *Xalan-Java is an XSLT processor for transforming XML documents into HTML, text, or other XML document types* »¹²⁹.

- Avec SDX-2, les sources de données peuvent être multiples, conduisant au concept générique d'entrepôt sans imposer une base de données extérieures au serveur.

La base de données et l'entrepôt sont deux plates-formes qui doivent être présentées à l'intérieur de SDX.

4.1.4. La base de données dans SDX

4.1.4.1. La base de données une définition

« La base de documents définit donc un ensemble de documents qui peuvent être recherchés. Cet ensemble partage des caractéristiques communes, en particulier un certain nombre de champs et leurs caractéristiques. On pourrait donc dire qu'une base de documents est définie fondamentalement par une liste de champs et par le contenu de ces différents champs, produit lors de l'indexation de documents »¹³⁰.

Un des avantages de SDX est qu'il n'impose pas de restriction sur les documents qu'il indexe au sein d'une base. Par contre il est possible d'indexer des documents XML respectant différents schémas ou DTD au sein d'une même base.

4.1.5. Le SDX et le multilinguisme

SDX accorde une grande importance à la langue et c'est pourquoi il est présenté comme un outil multilingue. « Autrement dit SDX est une plate-forme qui permet de réaliser des applications de recherche multilingues, cela ne signifie pas que SDX prévoit déjà tous les

¹²⁹ <http://xml.apache.org/xalan-j/>

¹³⁰ Sévigny, Martin et Frédéric Glorieux. Les bases de documents, <http://www.nongnu.org/sdx/docs/html/doc-sdx2/fr/presentation/bases.html> (dernier mise à jour, 2002/09/25)

outils nécessaires pour traiter correctement toutes les langues souhaitées par un développeur »¹³¹.

La langue dans SDX est construite autour de deux technologies (Java et XML) qui chacune a sa propre conception. C'est une « combinaison d'un code de langue ISO-639, d'un code de pays ISO-3166 (optionnel) et d'une variante (optionnelle). Pour représenter ces langues dans un contexte XML, l'attribut XML:lang est utilisé pour les deux premières parties et un attribut "variant" est utilisé pour le troisième. Exemple :

```
<sdx:firld xml:lang="fr-CA"variant="ac".../>1
```

Mais pour la langue d'interface, le code de pays est identifié à l'aide de l'attribut xml:lang. Pour construire les pages web dynamiques le SDX utilise en général le langage XSP dans lequel on trouve un élément qui s'appelle <sdx:page> qui permet d'utiliser par la suite l'un ou l'autre des services SDX. Exemple¹³² :

```
<sdx :document xml :lang="fr-CA">  
.....  
<sdx:results...>  
.....  
</sdx:results>  
</sdx:document>
```

L'exemple ci-dessus montre que le document XML aura, à sa racine, un élément <sdx:document> et un attribut xml:lang qui identifie la langue souhaitée pour l'interface. Cet attribut est toujours présent pour créer des interfaces multilingues.

Pour la langue de l'utilisateur et dans l'interface de gestion des utilisateurs de SDX, il est possible d'associer une langue à l'utilisateur. «Les applications qui ont leur propre interface de gestion des utilisateurs peuvent, de leur côté, associer une langue avec la méthode de leur choix, et SDX pourra tout de même utiliser cette langue pour indiquer la langue d'interface. »¹³³

¹³¹ Sévigny, Martin et Frédéric Glorieux. SDX et le multilinguisme. <http://www.nongnu.org/sdx/docs/html/doc-sdx2/fr/presentation/multilinguisme.html> (dernier mise à jour, 2002/09/25)

¹³² Sévigny, Martin et Frédéric Glorieux. SDX et le multilinguisme. <http://www.nongnu.org/sdx/docs/html/doc-sdx2/fr/presentation/multilinguisme.html> (dernier mise à jour, 2002/09/25)

¹³³ Sévigny, Martin et Frédéric Glorieux. SDX et le multilinguisme <http://www.nongnu.org/sdx/docs/html/doc-sdx2/fr/presentation/multilinguisme.html> (dernier mise à jour, 2002/09/25)

4.1.6. L'analyseur de mot

L'objectif d'un analyseur de mot est de transformer le contenu textuel dans le champ de type *word*. Après l'analyse, chaque mot individuel est stocké dans l'index pour distinguer les mots minuscules et pour supprimer les diacritiques. Mais il est difficile de trouver un analyseur de mot qui fonctionne bien pour toutes les langues car cela dépend de la langue de texte. Pour cette raison, le SDX permet au développeur d'application d'utiliser des analyseurs de mots différents pour tous les champs des bases de documents. «Ces analyseurs peuvent être soit inclus dans SDX, soit fournis par le développeur lui-même ; un analyseur de mots est tout simplement une classe Java qui étend (directement ou indirectement) la classe `fr.gouv.culture.sdx.search.lucene.analysis.Analyzer`.»¹³⁴.

Jusqu'à présent, SDX fournit des analyseurs pour cinq langues seulement, à savoir : la langue *chinoise*, la langue *allemande*, la langue *anglaise*, la langue *française* et la langue *russe*. Mais il a la possibilité d'accepter d'autres contributions pour d'autres langues car en structure interne, il fournit un mécanisme de configuration des analyseurs. Le fichier de configuration pour chaque langue se trouve sous l'adresse: `sdx/resources/conf/analysis` exemple¹³⁵ :

```
< ?xml version= "1.0" encoding="ISO-8859-1" ?>
<french useStopWords="true" keepAccents="false">
<stopWord>
<stopWord>le</stopWord>
<stopWord>la</stopWord>
<stopWord>les</stopWord>
</stopWord>
</french>
```

La langue arabe est un exemple du manque d'outils nécessaires pour la traiter au sein de SDX. Donc la base de donnée créée pour les manuscrits arabes sur SDX ne fournit pas une recherche par la langue arabe, espérons avoir le moyen d'intégrer les outils nécessaires pour effectuer cette tâche le plus tôt possible.

¹³⁴ Sévigny, Martin et Frédéric Glorieux. SDX et le multilinguisme. <http://www.nongnu.org/sdx/docs/html/doc-sdx2/fr/presentation/multilinguisme.html> (dernier mise à jour, 2002/09/25)

¹³⁵ Sévigny, Martin et Frédéric Glorieux. SDX et le multilinguisme. <http://www.nongnu.org/sdx/docs/html/doc-sdx2/fr/presentation/multilinguisme.html> (dernier mise à jour, 2002/09/25)

4.2. La base de données actuelle

Notre intention, dans cette partie, est de décrire la base de données hébergée sur le serveur de L'enssib et l'interface d'accès aux manuscrits arabes. Le but principal de cette base est de permettre aux utilisateurs d'effectuer des recherches parmi l'ensemble des champs sélectionnés. Une quarantaine de champs ont ainsi été sélectionnés parmi le 173 définis dans la DTD. « Le site devait donc permettre, pour le résultat d'une recherche, l'affichage structuré de documents descriptifs stockés dans un format sans mise en forme (XML) et l'affichage des images numérisées décrites dans le document XML »¹³⁶. La création de site Web a vu le jour grâce à un travail coopératif avec Guillaume Bourgois un étudiant de l'IUT Informatique de l'université Lyon 1, lors de son stage effectué à l'enssib en vue d'être utilisé plus tard à une plus grande échelle. Le page d'accueil si dessous montre les deux façons de recherche offertes par la base: Une recherche simple et une recherche avancée.



¹³⁶ Bourgois, Guillaume. Création d'une interface web pour l'indexation, la recherche et la consultation de manuscrits arabes anciens numérisés- Rapport de stage effectué à l'enssib du avril au juin 2003. 37pages. (p.9)

La recherche simple : elle ouvre une fenêtre cohérente à une série de sept petites icônes.

Par ces icônes il est ainsi possible de consulter:

- ❑ La liste des titres des manuscrits disponibles,
- ❑ La liste des auteurs dont un ou plusieurs manuscrits ont été indexés dans la base,
- ❑ La liste des copistes qui ont recopié, à partir d'un manuscrit original, un ou plusieurs manuscrits indexés dans la base de données,
- ❑ La liste des sujets traités dans les manuscrits,
- ❑ La liste des styles d'écriture que l'on peut trouver dans les manuscrits,
- ❑ La liste des pays où sont stockés les manuscrits papier,
- ❑ La liste des institutions qui conservent les manuscrits papiers.

En cliquant sur une de ces petites icônes, une liste fait apparaître sur l'écran les informations correspondantes (exemple : la liste des auteurs). Chaque élément apparaissant dans une liste contient un lien vers l'affichage d'un des documents si l'élément ne fait référence qu'à un seul document, comme par exemple un titre ou vers l'affichage de plusieurs documents dans le cas où l'élément désigne plusieurs manuscrits, comme dans le cas des listes des styles d'écriture. La figure ci-dessous montre l'écran avec les sept petites icônes en haut et un peu plus bas la liste des titres des manuscrits trouvés dans la base. Il est nécessaire de signaler ici que chaque fois que la base sera alimentée par des nouvelles informations, l'affichage de la liste correspondante sera mise à jour.

The screenshot shows the ENSIB website interface. At the top, the logo 'enssib' is displayed with the full name 'école nationale supérieure des sciences de l'information et des bibliothèques' underneath. Below the logo is a navigation bar with buttons for 'Accueil', 'Catalogue', and 'Recherche avancée'. To the right of these buttons is a search input field and a 'Rechercher' button. Below the navigation bar is a row of seven icons representing different search filters: 'Liste des titres', 'Liste des auteurs', 'Liste des copistes', 'Liste des sujets', 'Styles d'écriture', 'Liste des pays', and 'Liste des institutions'. The 'Liste des titres' icon is selected, and the corresponding list is displayed below. The list shows 10 items, each with a number and a title in French. At the top of the list, there is a dropdown menu set to '50' and the text 'résultats par page (n° 1 à 20 / 20)'. The list items are:

- 1 1ere volume de tafsir al-Qurāan al-Azim
- 2 al-Atar al-baqiya an al-qurun al-haliya.
- 3 al-Foutohat al-makiyeh
- 4 al-Qawaid al-fiqhiyeh
- 5 al-Taysir bisharh al-Jami el-saghir
- 6 Carte géographique portuaire de la mer méditerranée
- 7 Dala'el al khayrat wa shwariq al-anwaar fi dīkr al-salwat ala al-nabi al-mukhtar.
- 8 Diplôme confiant l'administration de l'église de sainte Mercure à Yuhanna ibn Gurgis
- 9 Discours sur les opérations manuelles : traité de chirurgie en trois parties
- 10 Kitab turyaq al-uqul fi'ilm al-usul,

Figure n°.62 : L'interface de la recherche simple dans la base de données des manuscrits arabes

La recherche avancée : la recherche avancée se compose de plusieurs champs qui aide les chercheurs à effectuer une requête composée de plusieurs éléments.

Advanced Search :

[Title]

Main Title

Incipit

Explicit

Volume Title

[From]

City

Region

Country

Date ([Hégire]) : from to

[Other]

Main Subject

Language

KeyWords

[Author]

Name

Profession

Born from to ([Hégire])

[Copyist]

Name

Born from to ([Hégire])

[Place of Conservation]

City

Region

Country

Identifier

Collection

Institution

Rechercher

Identification
Langue actuelle : English / Changer de langue : Français

Les étapes de cette page : 1) CSS 2) CSS 3) HTML

Figure n°.63 : L'interface de la recherche avancée dans la base de données des manuscrits arabe

La figure ci-dessus montre les six grands éléments de recherche : par titre, auteur, copiste provenance, lieu de conservation, et autre. Dans « **Titre** » il y a la possibilité d'effectuer une recherche par le titre principale, l'incipit, l'explicite et le titre de volume. Dans « **Auteur** », il est possible de chercher par le nom d'auteur, sa profession et sa date de naissance et/ou la date de sa mort soit par rapport à l'hégire, soit par rapport à la naissance de JC. Dans « **Copiste** » on peut chercher par son nom et par la date. La recherche par ville, région, pays et date se fait par l'aide d'éléments trouvés dans « **Provenance** ». Il est possible aussi de chercher par le lieu de conservation, (ville,

région, pays, identifiant, collection, établissement). Dans le dernier grand élément « **Autre** » il y a la possibilité de chercher par le sujet, langue et mots clés.

Il est possible aussi d'utiliser la liste déroutante pour choisir le sujet ou la date de publication par exemple. Bien évidemment, une recherche libre est aussi possible sur ces champs (sujet, date, pays, langue etc.).

4.2.1. *L'affichage de document*

Une fois qu'on a des résultats de recherche, l'affichage de ceux-ci sera donc possible. Sur le format d'affichage, on trouve seulement les champs décrivant un manuscrit donné. Par contre, pour certains champs on trouve un lien avec une URL (adresse de la page dans la base d'image). Ceci arrive où moment où l'image dans la base contient des caractéristiques qui correspondent à ce champ (exemple pour le champ illustration on va avoir une liste des URL pour toutes les images qui contiennent des illustrations dans le manuscrit). En cliquant sur cette adresse, une image sera donc affichée. Il faut signaler ici que tout ce travail de lien a été fait manuellement. Pour automatiser l'indexation, un travail avec l'INSA de Lyon est en cours de réalisation, il sera développé dans le chapitre suivant.

4.2.2. *L'affichage de base en multilingue*

Jusqu'à présent, la base de données assure un affichage multilingue anglais/français pour un maximum d'éléments. Un affichage en langue arabe serait idéal, mais à cause des problèmes liés à l'application de cette langue dans SDX, ce n'est pas encore réalisable. Un tel travail peut être possible dans le futur proche. Le site étant destiné à un accès mondial (les spécialistes de tout pays doivent pouvoir le consulter), la gestion d'un affichage multilingue est évidemment nécessaire.

4.3. Conclusion

SDX est un logiciel utile pour créer des documents XML sur l'internet. Sa flexibilité lui permet d'être très pratique pour le programmeur de créer leur propre base de données sans être strictement lié avec un logiciel pré défini qui se trouve sur le marché. Autre avantage de SDX est le fait qu'il soit gratuit. Cet avantage donne beaucoup d'espoir

pour les responsables des bibliothèques ou des centres de recherche de mieux gérer leur bibliothèque avec des dépenses très modestes. Notre base de données pour les manuscrits arabes est consultable par l'Internet grâce à ce logiciel