

D.E.S.S.
Conception Logicielle
et Applications aux Bases de Données

Grade de MASTER

Rapport de Stage

Portage de la plateforme logicielle SOX
et Création d'une Base de Données d'images

Effectué à :

Laboratoire Informatique, Image, Interaction (L3i)
Université de La Rochelle
Pôle Sciences et Technologie
17042 La Rochelle Cedex 1 – France

Sous la responsabilité de : Jean Marc OGIER – Professeur
Présenté par : Wilfrid PAPIN



Sommaire

Sommaire	1
Résumé	4
Introduction	5
I- Présentation du laboratoire L3i.....	6
A- Renseignements	6
B- Présentation	6
C- Les activités.....	7
D- Le L3i en quelques chiffres.....	8
1- Le personnel.....	8
2- Le budget	8
3- La production	9
E- Les partenaires.....	9
F- En savoir plus	9
II- Présentation du projet M.A.D.O.N.N.E	10
A- Objectifs et Contexte.....	10
1- Introduction.....	10
2- Enjeux du projet.....	11
B- Descriptif du projet.....	12
III- Positionnement du stage	13
IV- Description du stage	13
A- Contexte	13
B- Problématique.....	15
C- Objectifs techniques	17
V- Gestion de projet	18
A- Planification des tâches	18
B- Risque du projet	20
C- Consortium et réunions	21
VI- Environnement technique	22
A- Architecture logicielle.....	22
B- Données existantes	23
VII- Conception.....	24
A- Architecture du stage	24
B- Etude et fonctionnement de Docmining / SOX.....	25
1- Plateforme SOX : Fonctionnement.....	25
2- Problème majeur sur la plateforme SOX.....	26
3- Plateforme Docmining : Fonctionnement.....	26
C- Outil de création pour l'interface	29
D- Mécanisme de portage.....	30
E- Moyen de stockage des images	33
VIII- Réalisation.....	37
A- Développement du module Interface	37
B- Intégration du module de traitement à l'interface	42
C- Portage de la plateforme logicielle sur Internet.....	42
D- Création du moyen de stockage	44
E- Création du mini site	45
1- Recherche d'image par critères.....	45

a) Par signature	45
b) Par méta-données	48
2- Affichage Répertoire.....	50
3- Insertion dans la base	50
IX- Objectifs fixés et résultats obtenus	51
X- Acquis techniques	53
XI- Bilan humain.....	54
Conclusion.....	55
Table des légendes	56

Résumé

Dans le but de la sauvegarde et de la valorisation des données patrimoniales, certains membres du laboratoire L3I (Laboratoire Informatique, Image, Interaction) ainsi que des membres de laboratoires partenaires ont répondu à un appel à projet du Ministère de la Recherche et du CNRS sous la forme d'une ACI et ont monté le projet M.A.D.O.N.N.E.

Ce projet s'adresse au patrimoine documentaire papier et il vise à mettre en place des outils automatiques d'analyse et d'indexation des images de documents pour faciliter la navigation dans les bases documentaires habituellement inaccessibles sous forme papier.

Dans ce cadre, les partenaires du consortium M.A.D.O.N.N.E ont été amenés à développer une plateforme logicielle, nommée SOX, permettant d'intégrer des outils hétérogènes de traitement. Afin de faciliter les travaux de tous les partenaires, l'accès direct de la plateforme sur Internet serait un outil de travail collaboratif essentiel pour la promotion des travaux du consortium.

Afin de réaliser ceci, le consortium a défini ses besoins et ses souhaits en terme d'utilisation de la plateforme. Celle-ci est présentée sur le site MADONNE et permet d'exécuter divers traitements sur un document que l'utilisateur de la plateforme aura choisi au préalable.

Le laboratoire L3i fait partie de l'université de La Rochelle dont j'ai été un des étudiants il y a un an. Dans ce cadre, j'ai postulé pour un stage en m'adressant à Jean Marc Ogier. Il me proposa à travers ce projet MADONNE un sujet de stage correspondant à ces besoins. La diversité des compétences nécessaires à la réalisation de ce projet m'a poussé à le choisir. L'analyse de besoins et des souhaits m'a conduit au développement de la plateforme et de son portage sur Internet. Et afin de gérer de manière rapide et efficace les images utilisées lors de l'exécution de ses traitements, une base de données MYSQL a été créée permettant de gérer ces images via différents critères.

Finalement, le travail réalisé à travers ces deux modules doit permettre de promouvoir les travaux du consortium afin de poursuivre dans le bon sens les recherches sur la valorisation du patrimoine français.

Introduction

Le stage se déroule au sein du Laboratoire Informatique, Image, Interaction (L3i) de l'Université de La Rochelle. Au quotidien, le laboratoire travaille sur des grands projets dans le domaine de l'Image et du Comportement.

Dans ce cadre, le projet M.A.D.O.N.N.E a vu le jour via une Action Concertée Incitative (ACI). Cette ACI émane d'un regroupement des laboratoires français spécialisés dans le domaine de l'analyse des images de documents comme le Laboratoire L3i de La Rochelle ou le Laboratoire LORIA à Nancy. L'objectif de cette ACI est de proposer une réflexion prospective, un cadre méthodologique et des techniques de structuration des contenus des images pour permettre d'organiser et de préparer le déploiement à grande échelle de cette activité de numérisation et dans le but d'atteindre une numérisation de qualité suffisante pour la création de bases partageables de contenus. Pour répondre à ces besoins, le consortium MADONNE a mis au point il y a quelques temps une plateforme logicielle, nommée SOX, permettant d'intégrer des outils hétérogènes de traitement. Ces traitements s'appliquant sur des documents de formats différents tels que PDF, PostScript, etc.

A mon arrivée, l'utilisation de cette plateforme était cependant limitée car son utilisation ne pouvait se faire qu'en local. Par conséquent, un laboratoire partenaire du consortium pouvait travailler sur un nouveau traitement et l'insérer dans la plateforme, les autres laboratoires ne puissent en profiter automatiquement et rapidement. Dans cette optique, l'ACI M.A.D.O.N.N.E a choisi de porter cette plateforme logicielle sur Internet. Cela permet entre autres à la plateforme d'être mise à jour plus facilement et donc facilite son utilisation.

En parallèle à ce besoin, Jean Marc Ogier souhaite gérer les images que le consortium utilise de manière plus efficace. Pour le moment, les images se trouvent pour la plupart sur le serveur FTP du CESR (Centre d'Etudes Supérieures de la Renaissance de Tours). Ce centre fournit toutes les images de documents anciens selon les besoins du consortium. Pour ce faire, il souhaite créer une base de données d'images avec les méta-données s'y rapportant. Sous forme de stage, ce développement de deux modules m'a été proposé.

Dans ce cadre, ce stage a pour but de concevoir, réaliser et mettre en place ces modules. Concernant la création de la base, la solution s'appuie sur une gestion des chemins d'accès aux images dans la base et non à l'image en elle-même. Et pour ce qui est de la plateforme logicielle, la solution se tourne vers la création d'une applet.

La partie suivante de ce rapport présente le laboratoire L3i. Nous abordons ensuite la description générale du projet M.A.D.O.N.N.E puis celle du stage avant de détailler les aspects de la gestion du projet. La suite du rapport devient plus technique puisqu'elle introduit l'environnement matériel et logiciel utilisé et les différents aspects de la phase de conception. Une dernière partie est enfin consacrée à la réalisation en elle-même, jusqu'à la mise en production finale.

I- Présentation du laboratoire L3i

A- Renseignements

Le laboratoire L3i (Laboratoire Informatique, Image, Interaction) est reconnu par le Ministère de la Recherche comme Équipe d'Accueil (EA2118) depuis 1997.



Ses coordonnées sont les suivantes :

L3i Université de La Rochelle
Pôle Sciences et Technologie
17042 La Rochelle Cedex 1 - FRANCE
Tél : 05 46 45 82 62
Fax : 05 46 45 82 42
Tel international : +33 5 46 45 82 62
Site Internet : <http://www-l3i.univ-lr.fr/>

B- Présentation

Le laboratoire L3i regroupe désormais une quarantaine d'enseignant chercheurs (dont 12 habilités à diriger des recherches), essentiellement issus de la communauté informatique (section 27 du CNU) et génie informatique (section 61) de l'Université de La Rochelle. Au premier janvier 2004, 20 doctorants sont inscrits.

Le laboratoire a changé de nom il y a quelques années. De Laboratoire d'«*Informatique et Imagerie Industrielle*», il s'est tourné vers «*Informatique, Image, Interaction*». Il se positionne donc beaucoup plus sur une image, facteur d'interactions que sur les aspects technologiques associés à une informatique Industrielle. Cette orientation a été motivée par une volonté d'aborder, de la manière la plus cohérente, complète et fondamentale, les différentes facettes de chaque contexte applicatif auquel le laboratoire est confronté. En particulier, l'*ouverture du laboratoire aux autres domaines disciplinaires* (les arts, le littoral, les usages) a montré que la prise en compte des interactions, sous toutes les formes qu'elles prennent, était primordiale.

C- Les activités

L'Image et le Comportement sont donc au centre de l'activité du Laboratoire.

Une partie des activités du laboratoire concerne donc naturellement les aspects analyse/modélisation/synthèse – c'est l'objet de la thématique *Image et séquences d'Images (ISI)*. Cette thématique s'est tout d'abord focalisée sur la problématique du traitement des images acquises à partir de systèmes d'imagerie de type industriel. Cette problématique s'est donc déclinée très rapidement sur le mode restauration, extraction d'informations pertinentes et détection selon les spécificités industrielles: conditions d'acquisition difficiles, contraintes d'éclairage, contraintes de temps réel, image de haute définition. Elle s'est par la suite enrichie de réflexions portant sur des domaines d'imageries *multi-composantes* (couleur, multi-canaux, multi-modalités, multi-formes) et s'appuyant sur des séquences temporelles et/ou sur des contextes de systèmes d'information documentaires.

Pour l'analyse d'Image et pour l'interprétation de comportements complexes, un travail de fusion/classification est nécessaire. La thématique *Données, Formes, Interprétation (DoFin)* traite ce problème et le replace dans un cadre plus général, lui permettant d'aborder, en particulier, les données semi structurées. L'objectif de cette thématique est la conception de modèles robustes et le développement d'outils pour l'analyse de données et la structuration d'informations complexes. Du numérique au symbolique, différents niveaux sémantiques données statiques/dynamiques, spatiales/temporelles, relations spatio-temporelles, informations contextuelles, etc.), différents niveaux de structuration (indices visuels basiques, trajectoires, données hétérogènes, etc.) sont pris en compte dans la modélisation, ainsi que la gestion de l'imperfection (incertitude, imprécision, incomplétude). Pour cela, DoFin est organisée autour de trois thèmes :

- Reconnaissance des Formes ;
- Analyse et Interprétation de Données Spatio-temporelles ;
- Structuration pour la Recherche d'Information, et le contexte applicatif concerne essentiellement l'image, la séquence d'images et les données environnementales ;

La thématique *Modèles, Comportements, Architectures (MoCA)* contribue à inscrire les travaux menés dans un contexte opérationnel et un cadre méthodologique. Ces objectifs sont donc la conception de modèles, de méthodes, de langages et d'outils pour représenter, analyser et mettre en œuvre le comportement dans des systèmes complexes, mobiles et évolutifs. Les activités de recherche sont organisées en trois classes de préoccupations complémentaires :

- Démarche de conception : Travaux sur la cohérence et l'optimisation, recherche de critères de qualité et de métriques ;
- Modélisation et conception : Intégration de modèles, modèles de comportement, définition de composants, validation de composants, automatisation de production de code ou de modèles intermédiaires, modélisation à base de systèmes multi agents, modèles de données pour les systèmes de visualisation d'informations spatiales embarquées ;
- Production d'outils théoriques et cadre méthodologique : treillis de Galois et ordres partiel, systèmes implicatifs ;

Ceci permet de mieux comprendre, à travers leur modélisation, les systèmes représentés par des images, mais aussi, de produire les logiciels associés à nos résultats en les intégrant dans une architecture pertinente.

D- Le L3i en quelques chiffres

1- Le personnel

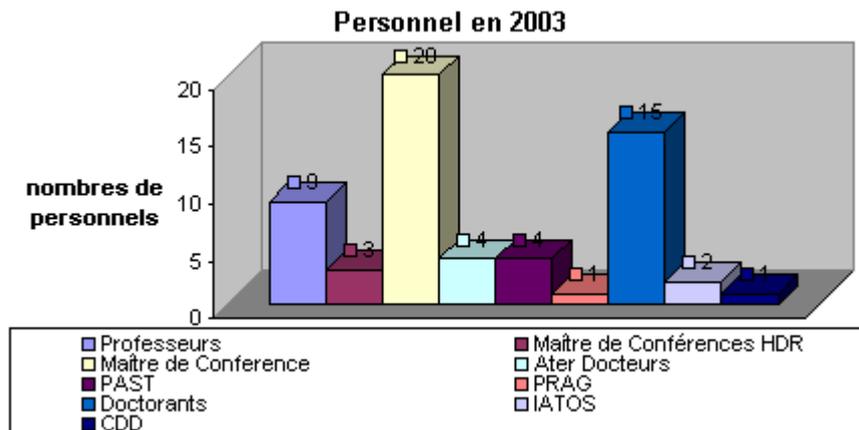


Figure 1 - Le personnel du L3i en 2003

2- Le budget

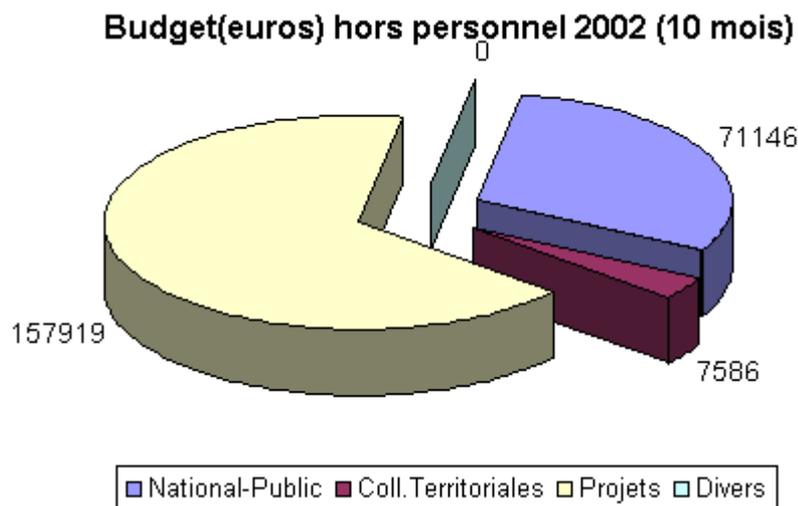


Figure 2 - Le budget du L3i en 2002

3- La production

Les publications du laboratoire ont gagné en qualité et en quantité; ceci est, en particulier, observable sur la production de cette dernière année. Une vraie dynamique de publication a été créée.

Publications/années	1999	2000	2001	2002 (10 mois)
Revue Internationales	3	7	7	11+9
Communications actes	19	26	20	21
Conférences Invitées	0	0	1	0
Autres	0	3	3	2
Thèses+HdR	1	3	5+1	2+2

E- Les partenaires

Le laboratoire axe ces recherches sur des travaux en liaison directe avec les acteurs régionaux (industriels, communauté de Villes, pôle Image d'Angoulême, Futuroscope, ...) ou encore des partenaires universitaires ou industriels nationaux ou internationaux (France Télécom R&D, Université de Rouen, Université de Fribourg, Université en Malaise, ...).

F- En savoir plus

- Le laboratoire : <http://www-l3i.univ-lr.fr/>
- L'Université de La Rochelle : <http://www.univ-lr.fr/>

II- Présentation du projet M.A.D.O.N.N.E

A- Objectifs et Contexte

1- Introduction

De nombreux débats d'experts européens ont récemment émis de fortes recommandations concernant les actions à entreprendre pour favoriser de manière durable la coordination et la valorisation des activités de numérisation du patrimoine. En effet, les ressources culturelles et scientifiques de l'Europe sont un bien public unique qui représente la mémoire collective et vivante de nos différentes sociétés et qui forme une base solide pour le développement des industries au contenu numérique dans une société de la connaissance durable.

Sur la base de la réunion préparatoire de Luxembourg des 15 et 16 novembre 2000, ces experts ont souligné la valeur et l'importance du contenu numérisé européen à considérer dans les domaines culturel et scientifique comme :

- un patrimoine accessible et durable. L'Europe possède un héritage culturel et scientifique d'une richesse unique. La numérisation de ces ressources est fondamentale pour que les citoyens puissent y avoir davantage accès et que l'héritage culturel collectif de l'Europe (passé et à venir) soit préservé ;
- un atout en faveur de la diversité culturelle, de l'enseignement et des industries de contenu. La numérisation des biens culturels est essentielle au maintien et à la promotion de la diversité culturelle dans un contexte de mondialisation. Les données numérisées constituent également des ressources précieuses pour l'enseignement et les industries du tourisme et des médias ;
- un ensemble de ressources numérisées riches et variées.

Les États membres de l'Union Européenne ont beaucoup investi dans les programmes de numérisation du contenu culturel et scientifique. Ces opérations de numérisation portent à la fois sur une variété de domaines et sur une variété de types de contenus, tels que les collections des musées, les sites archéologiques, les archives audiovisuelles, les cartes, les documents historiques et les manuscrits.

Les programmes découlant de ces études viseront à créer une infrastructure pour la culture numérisée afin de permettre l'accès au patrimoine culturel et scientifique numérisé. Parmi les tâches identifiées, les objectifs sont :

- améliorer la qualité et la facilité d'utilisation du contenu,
- multiplier les possibilités d'accès offertes à la population,
- et renforcer les actions de sensibilisation sur les problèmes de conservation à long terme des données.

Ces premiers objectifs ne pourront être atteints que :

- en concluant des accords sur des normes d'interopérabilité,
- en définissant des lignes directrices pour la conservation des données et la pérennité de la description des contenus,
- en élaborant des modèles cohérents,
- en adoptant de bonnes pratiques pour la gestion des droits et des biens, ainsi qu'en développant des modèles commerciaux cohérents pour la culture numérisée.

La tâche est large et ne peut s'accomplir sans la volonté des pouvoirs publics. Et ce sont différentes instances académiques d'objectifs régionaux, et complémentaires, qui espèrent l'accomplir dans l'établissement d'une collaboration nationale française au minimum, et pour une reconnaissance scientifique internationale des activités de la recherche française dans ce domaine.

2- Enjeux du projet

Sur la base de ces éléments qui mettent en évidence la nécessité du déploiement de nouvelles mesures de numérisation cohérentes, les grandes orientations de l'ACI M.A.D.O.N.N.E (Action Concertée Incitative) pour la valorisation du patrimoine s'appuieront sur les notions suivantes.

Le coût d'une numérisation de qualité est tel qu'il faut pouvoir assurer la pérennité et la ré-utilisabilité des contenus numériques ainsi créés sur le long terme. Or actuellement aucune approche de numérisation n'intègre cette démarche, la numérisation ne se limitant trop souvent qu'à la simple capture en masse d'images. Les échecs de quelques expériences industrielles et institutionnelles démontrent d'autre part que des opérations de numérisation ne peuvent être décorréliées de projets intégrant tous les maillons de la chaîne de numérisation (de l'acquisition des images à la production de divers contenus numériques). L'une des propositions du consortium est donc la constitution et le partage, à un niveau qui se devra à terme de dépasser le simple cadre national, d'entrepôts d'images des patrimoines culturels européens. L'objectif de cette ACI est de proposer une réflexion prospective, un cadre méthodologique et des techniques de structuration des contenus des images pour permettre d'organiser et de préparer le déploiement à grande échelle, à long terme et pour le plus grand nombre de cette activité de numérisation qui tend à se généraliser, dans le but d'atteindre une numérisation de qualité suffisante pour la création de bases partageables de contenus. Il s'agit, sur la base de l'expertise des équipes françaises spécialisées dans la numérisation de documents, de proposer une approche fondatrice pour aider à contrôler cette activité. Cependant et bien que cette proposition déborde largement le seul cadre des documents patrimoniaux, nous limiterons notre projet au seul cas d'images de ce type (pourtant très large).

D'autre part, cette orientation de la numérisation du patrimoine nécessite de pouvoir anticiper sur le long terme les besoins et les usages des contenus patrimoniaux tout en assurant leur préservation numérique (archives numériques). Dans cette perspective, il s'agit de proposer une démarche capable, par exemple, de répondre à la demande, à un besoin spécifique d'intégration de données patrimoniales dans un Système d'Information (administratif, touristique, de bibliothèque ...).

Nous situons notre démarche en dehors du Système d'Informations car à ce niveau l'indexation est en général trop dépendante d'une activité spécifique. De même en amont de la chaîne de numérisation, la capture des données (scanneurs, caméras) n'est pas considérée car elle est soumise à une évolution des technologies importante pour la considérer pérenne. La capture est également très dépendante d'une organisation de service (archives, musées, collectivités locales) et donc soumise à des évolutions externes importantes par rapport à l'activité qui est visée à long terme.

Ce projet se focalise donc sur la partie concernant l'indexation, l'organisation et l'enrichissement progressif et incrémental des entrepôts de données patrimoniales, dans un but de production de services génériques aux usagers, sur la base de garanties d'interopérabilité entre données, mais également entre traitements. S'appuyant sur des collections de documents après acquisition numérique, la problématique est ici de déployer un

ensemble d'outils génériques et transversaux permettant d'annoter numériquement les documents de manière progressive, permettant ainsi à différents usagers de naviguer dans les bases documentaires. Les collections de documents hétérogènes se verront donc, de manière incrémentale, affectées des méta-données permettant la navigation intra et extra-collection, sur la base de mécanismes d'indexation croisée.

Pour finir, cette ACI émane d'un regroupement des laboratoires français spécialisés dans le domaine de l'analyse des images de documents. Ces structures font partie de réseaux organisés (journaux, conférences, GRCE, GDR I3, AS numérisation, ISDN de Lyon...) et ont plusieurs succès à leur actif : lecture des adresses postales, des formulaires, des plans cadastraux, des imprimés, reconnaissance des structures des documents, rétro-conversion XML des sommaires et catalogues, indexation automatique des images de documents. Par ailleurs citons leur expérience dans le domaine des documents anciens : projet DEBORA, archives de la Mayenne, traitement des manuscrits de Flaubert, ...

B- Descriptif du projet

Comme nous venons de l'évoquer, la valorisation des collections du patrimoine soulève de nombreuses difficultés essentiellement liées à la grande variabilité des contenus et des usages qui peuvent être constatés sur ces collections numérisées. La réflexion qui est proposée ne vise pas à apporter une réponse précise et définitive à l'ensemble des activités liées à ces fonds patrimoniaux (impossibilité d'anticiper l'ensemble des activités sur le long terme) mais à proposer, sur la base des modes de consultation constatables ou envisageables, des réponses pertinentes du point de vue des techniques de traitement d'images et de reconnaissance des formes, pour structurer les bases d'images de documents du patrimoine. Pour répondre à cet objectif, 5 sous-thèmes ont été identifiés comme déterminants. Chacun de ces 5 sous-thèmes est décliné suivant la nature des problématiques associées, les verrous scientifiques correspondants, et les objectifs visés pour passer outre ces contraintes actuelles.

- 1 - *Modélisation des collections*
- 2 - *Indexation sur les structures*
- 3 - *Indexation par le contenu*
- 4 - *Adéquation Modèle-Traitement*
- 5 - *Compression*

III- Positionnement du stage

Le projet MADONNE s'inscrit dans un domaine de recherche et a pour but de proposer des techniques de traitement d'images, de reconnaissance de formes. Le stage qui m'a été proposé sort totalement de ce contexte au point de vue recherche. Il s'agit de promouvoir les travaux effectués par chacun des partenaires du consortium afin de montrer les possibilités du consortium en matière de recherche et bien évidemment, afin de faire évoluer les techniques d'imagerie.

Dans ce cadre, le consortium MADONNE m'a proposé de créer des outils permettant de promouvoir leurs travaux. Donc ce stage s'inscrit dans une optique moins recherche et se positionne plus en périphérie du projet M.A.D.O.N.N.E.

IV- Description du stage

A- Contexte

À l'ère de la circulation massive d'informations de tous ordres sur les réseaux mondiaux, la pauvreté structurelle de toute l'information disponible devient un problème crucial. Par "pauvreté" de l'information, nous entendons l'absence de structuration permettant une indexation et une interprétation orientée navigation. Un premier exemple venant à l'esprit est, bien entendu, le document sous forme papier, scanné et stocké comme une image dans une archive. Mais d'autres formats de représentation de bas niveau (PostScript, PDF, DXF, voire des pages HTML classiques) posent les mêmes types de problèmes quand il s'agit d'organiser l'information utile, de l'indexer, de la retrouver aisément. Cette information est caractérisée par une composition de parties textuelles, d'illustrations graphiques et de photographies. Nous nous intéresserons principalement aux deux premières composantes, et ne chercherons pas à extraire des critères d'indexation par le contenu des images de type photographique (c'est un autre sujet de recherche à part entière).

Les verrous technologiques abordés sont ceux liés à l'accessibilité et à l'acceptabilité de l'information par un service. Ils reposent sur le fait qu'il n'existe pas actuellement de système fiable autorisant de faire un conversion de documents sans prendre en compte des connaissances métiers formalisées dans un modèle.

Dans ce cadre, les partenaires du consortium M.A.D.O.N.N.E ont été amenés à développer une plate-forme logicielle, nommée SOX, permettant d'intégrer des outils hétérogènes de traitement (traitement d'image, de reconnaissance des formes, d'intelligence artificielle), développés dans des environnements et langages hétérogènes (C, C++, Java, Windows, Java, ...).

Une première version du démonstrateur SOX (Segmentation Orientée Xml) a été réalisée. Cette plate-forme de traitement de documents comporte différentes briques logicielles :

- Des bibliothèques de traitements provenant des différents membres du consortium ou de logiciels libres.
- Un moteur d'interprétation de scénario de traitements.
- Des outils pour la définition de scénario.
- Des outils pour intégrer de nouveaux traitements à partir de logiciels externes
- Des outils pour intégrer de nouveaux traitements à partir de scénarios SOX
- Des outils pour la visualisation des résultats d'un scénario.
- Divers outils de manipulation de fichiers image et XML.

Le moteur de SOX est l'interprétation de scénarios. Un scénario décrit sous forme d'arbre les traitements à réaliser ainsi que l'information qu'ils manipulent.

Les documents traités peuvent se trouver sous forme de fichier image ou sous forme de fichier électronique incluant une partie d'image. Le but est d'obtenir en résultat de traitement d'un document par un scénario, sa représentation sous forme d'arbre d'objets au format XML.

La technologie XML a été choisie non seulement comme format de représentation du document mais aussi pour le paramétrage des traitements, des scénarios et des paramètres de configuration de l'application. SOX est développé entièrement en Java, même s'il intègre des bibliothèques de traitements écrites en C/C++ ou des traitements fournis sous forme d'exécutable.

L'objectif du premier sujet du stage consiste à porter cette plateforme SOX en servlet ou JSP (Java ServerPage) pour être utilisée à distance via un navigateur WEB.

La seconde partie du stage ne concerne pas cette plateforme logicielle SOX mais seulement les images qu'elle traite. En effet, les traitements d'image, de reconnaissance des formes, d'intelligence artificielle, intégrés dans la plateforme, sont testés sur des documents très différents allant du document PDF pour la reconnaissance de caractère à la simple image au format BMP, JPEG ou TIFF pour la reconnaissance des formes.

Tous les partenaires du consortium travaillent en phase expérimentale sur les mêmes images venant du Centre d'Etudes Supérieures de la Renaissance de Tours. Ce centre travaille sur les documents anciens tels que les manuscrits ou les documents religieux qu'ils prennent en photo par le biais d'appareil photo numérique de très grande qualité afin d'obtenir des images ayant une définition et une qualité exceptionnelle. Actuellement, le laboratoire L3i de La Rochelle stocke de nombreuses images du CESR sur leur serveur. Mais ces images ne sont ni cataloguées ni répertoriées et donc il est souvent très difficile de trouver des images précises ou même de comparer des images entre elles via des traitements.

Donc l'objectif de la deuxième partie de mon stage étant de créer une base de données de toutes ces images via l'outil le plus adapté, et de réaliser un mini site permettant la recherche d'images par le biais de critères que l'utilisateur choisira, principalement basés sur le contenu de l'image.

B- Problématique

Cette plateforme logicielle SOX permet d'interpréter des scénarios dans le but d'obtenir une représentation d'un document de type PDF par exemple sous forme d'objets au format XML.

Voici la plateforme SOX permettant de créer des scénarii et de les interpréter :

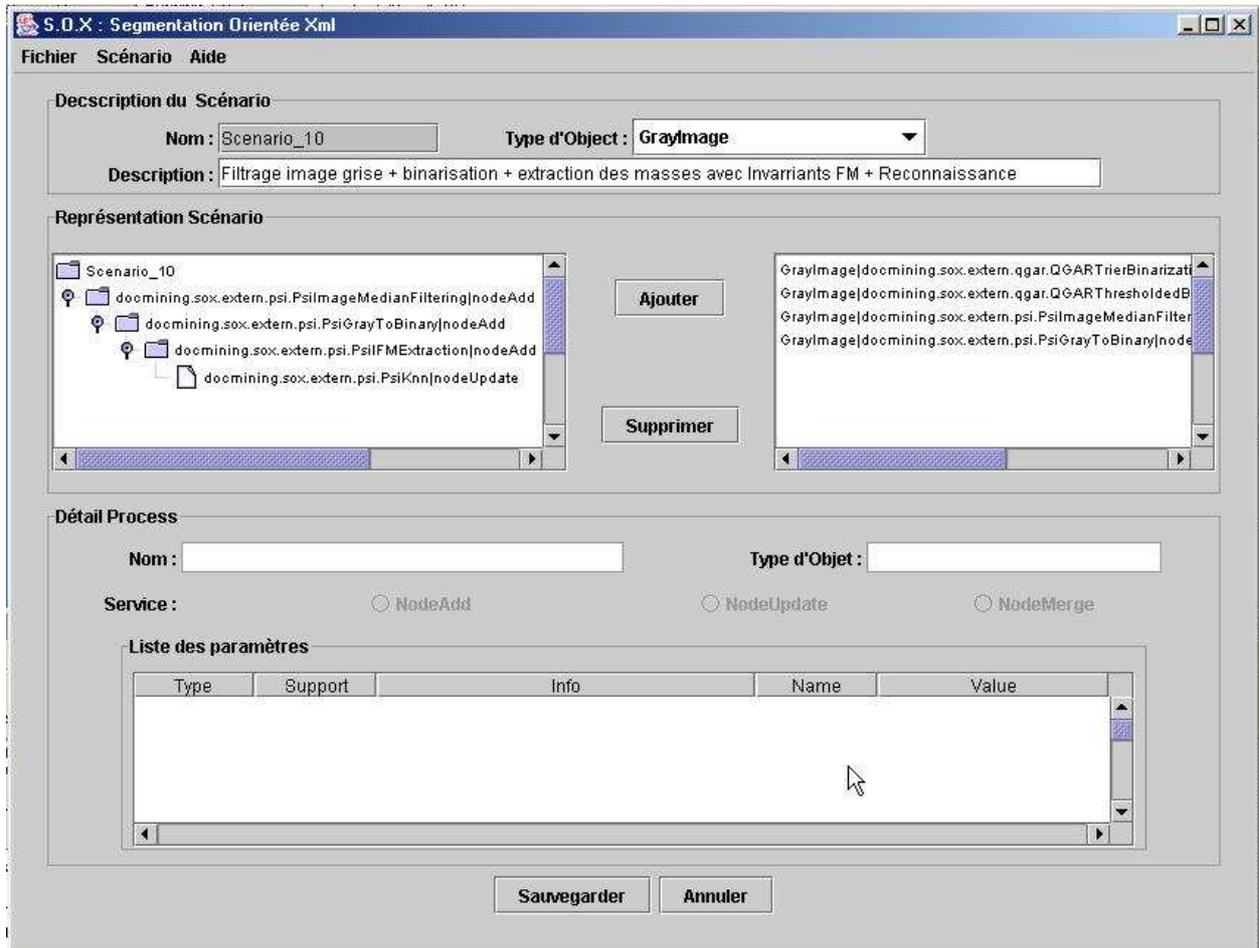


Figure 3 - Plateforme SOX

Cette plateforme contient des composants lourds comme :

- JFileChooser
- JTree
- JTable
- JList

Le résultat de l'interprétation du scénario par la plateforme SOX donne le résultat suivant :

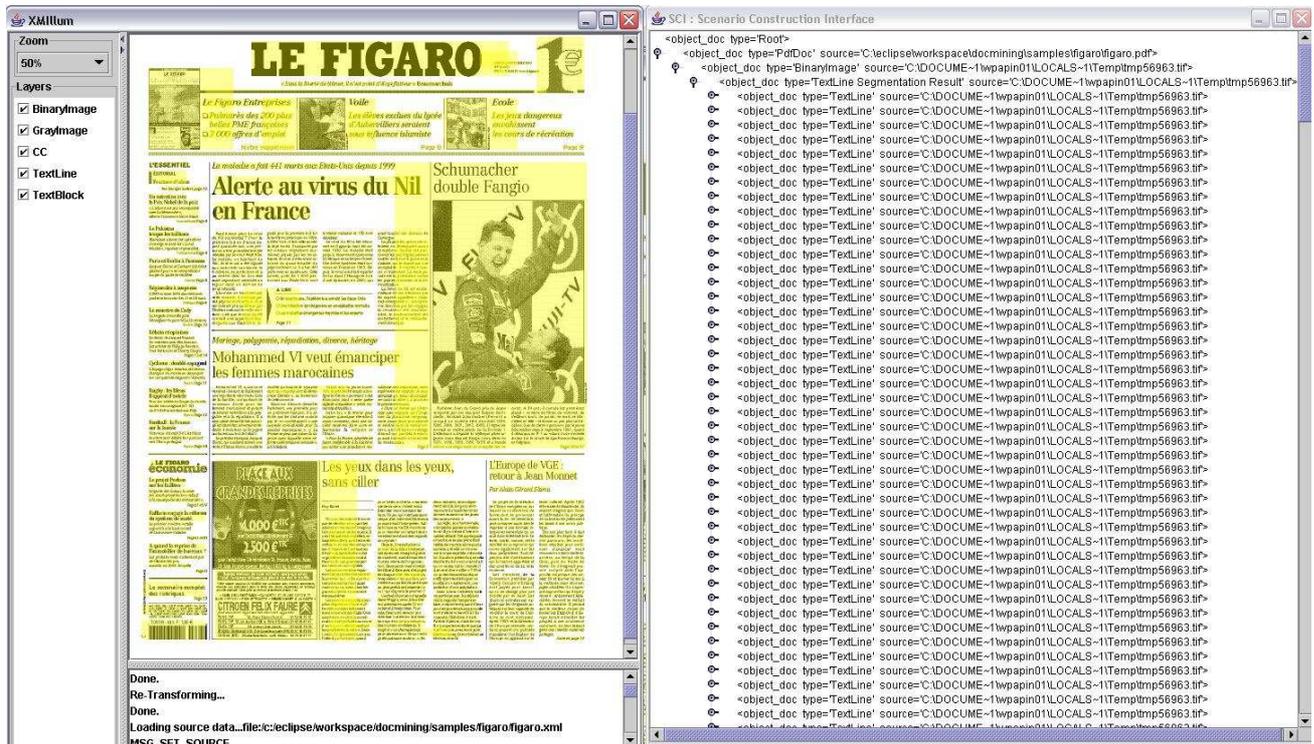


Figure 4 - Résultat de l'interprétation du logiciel SOX

Le résultat se présente sous la forme de deux fenêtres :

- l'une contient le document à traiter
- l'autre, le résultat du traitement sous forme XML (via un arbre)

L'utilisateur a la possibilité d'effectuer de nouveaux traitements via les deux fenêtres ci-dessus. Il lui suffit de sélectionner un objet dans l'arborescence XML et de faire un clic droit sur la souris pour faire apparaître les traitements que l'on peut lancer par rapport à l'objet sélectionné.

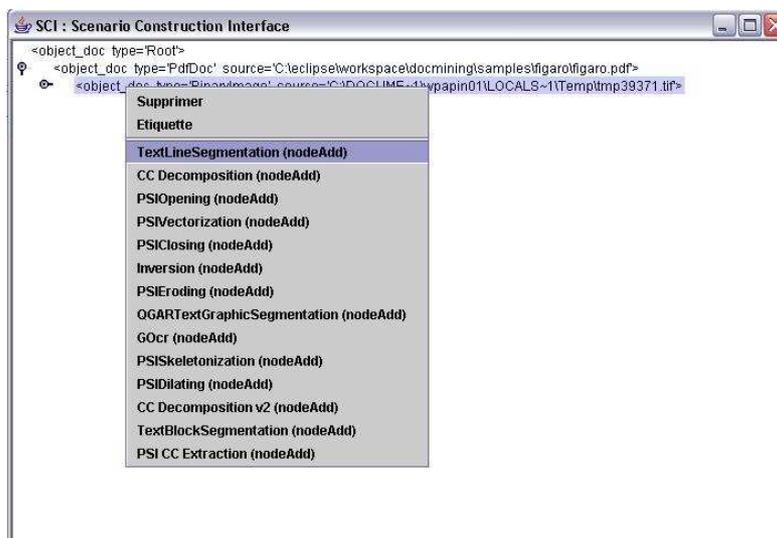


Figure 5 - Menu de traitements possibles sur un objet

Les besoins du Consortium M.A.D.O.N.N.E par rapport à cette plateforme sont :

- l'utilisation de la plateforme logicielle sur le Web ;
- l'intégration d'outils hétérogènes de traitement directement via la plateforme et sur le Web ;

Les besoins du consortium par rapport à la base de données :

- la gestion de toutes les images que gère le consortium via une base de données ;
- la recherche multicritères des images (critère tel que par la signature ou les méta-données);

Le portage de cette plateforme sur Internet a un double objectif :

- un aspect promotionnel, permettant de donner un aperçu au monde de la recherche de ce qu'il est possible de faire avec les nouvelles technologies ;
- un aspect technique permettant à tous les partenaires du Consortium de tester la plateforme directement sur Internet ;

La création de la Base de Données a également un objectif :

- un aspect temps de recherche, permettant de réduire de manière significative le temps de recherche d'images ;

La vision métier du projet exprime les besoins souhaités et les objectifs à atteindre. La partie suivante en exprime les aspects techniques.

C- Objectifs techniques

Dans ce projet, différents points sont à développer :

- recherche du meilleur outil pour le portage de la plateforme sur le Web ;
- création de la partie interface de la plateforme (car problème majeur sur la plateforme SOX, voir la partie Conception de ce rapport) ;
- création de la partie portage en utilisant l'outil choisi au préalable ;
- création de la base de données d'images ;
- développement des outils permettant de travailler sur la base de données ;

Les compétences techniques identifiées pour répondre à cet aspect sont :

- le développement du portage en Servlet ;
- le développement de l'interface de la plateforme en Java Swing ;
- une base de données MySQL pour le stockage et la gestion des images ;

V- Gestion de projet

Cette partie présente les aspects concernant la gestion de projet et l'organisation du travail. Un plan précis s'appuyant sur différentes phases a été suivi : la planification des tâches, l'étude et le suivi des risques, et l'organisation de réunions internes afin de veiller au bon avancement du projet.

A- Planification des tâches

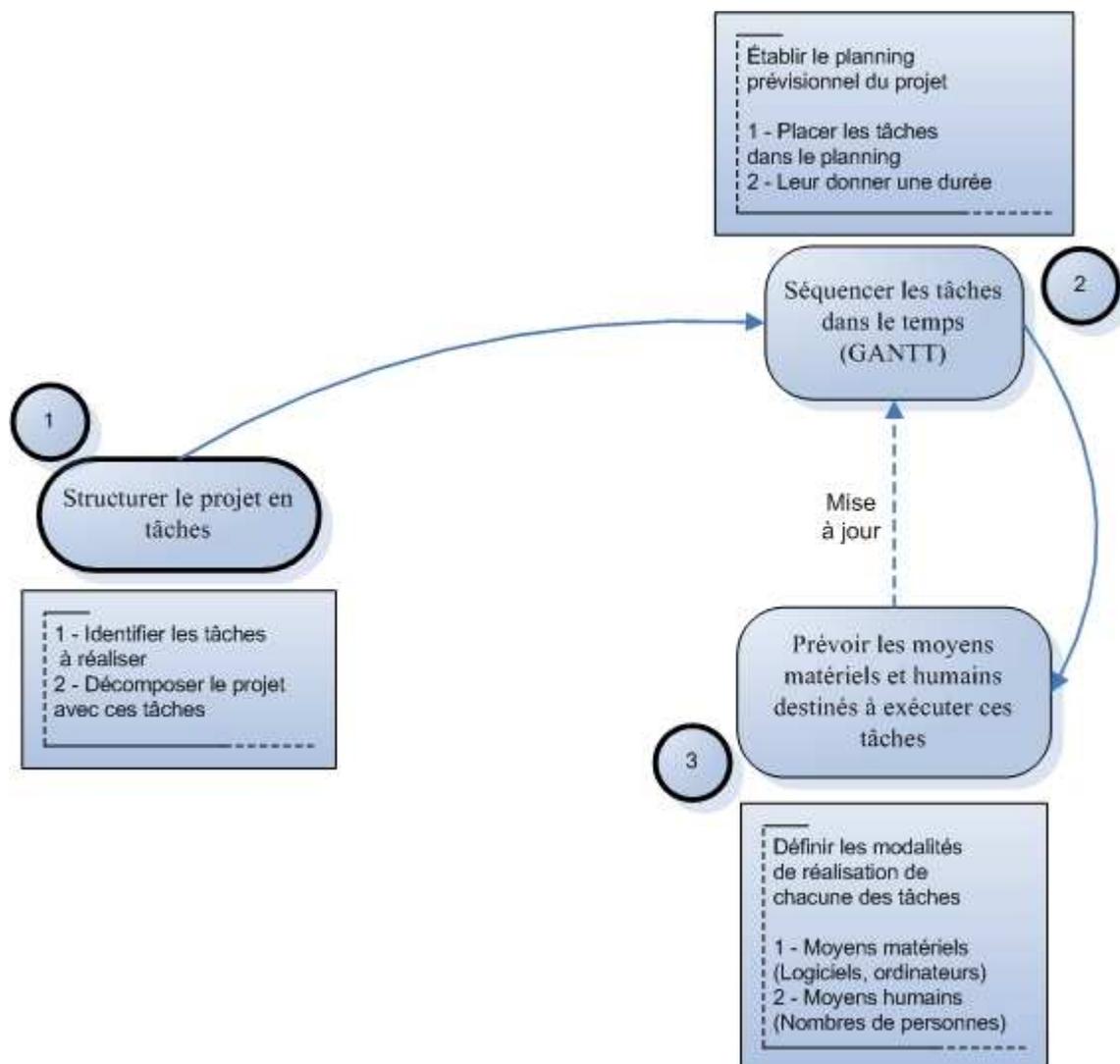


Figure 6 - Diagramme de planification

Cette démarche permet d'élaborer un planning des tâches représenté par le biais du logiciel Microsoft Project et via l'utilisation du diagramme de GANTT. Le diagramme de la figure 7 présente la liste des tâches établies pour la réalisation du stage mais également sa planification pour la durée de ce dernier.

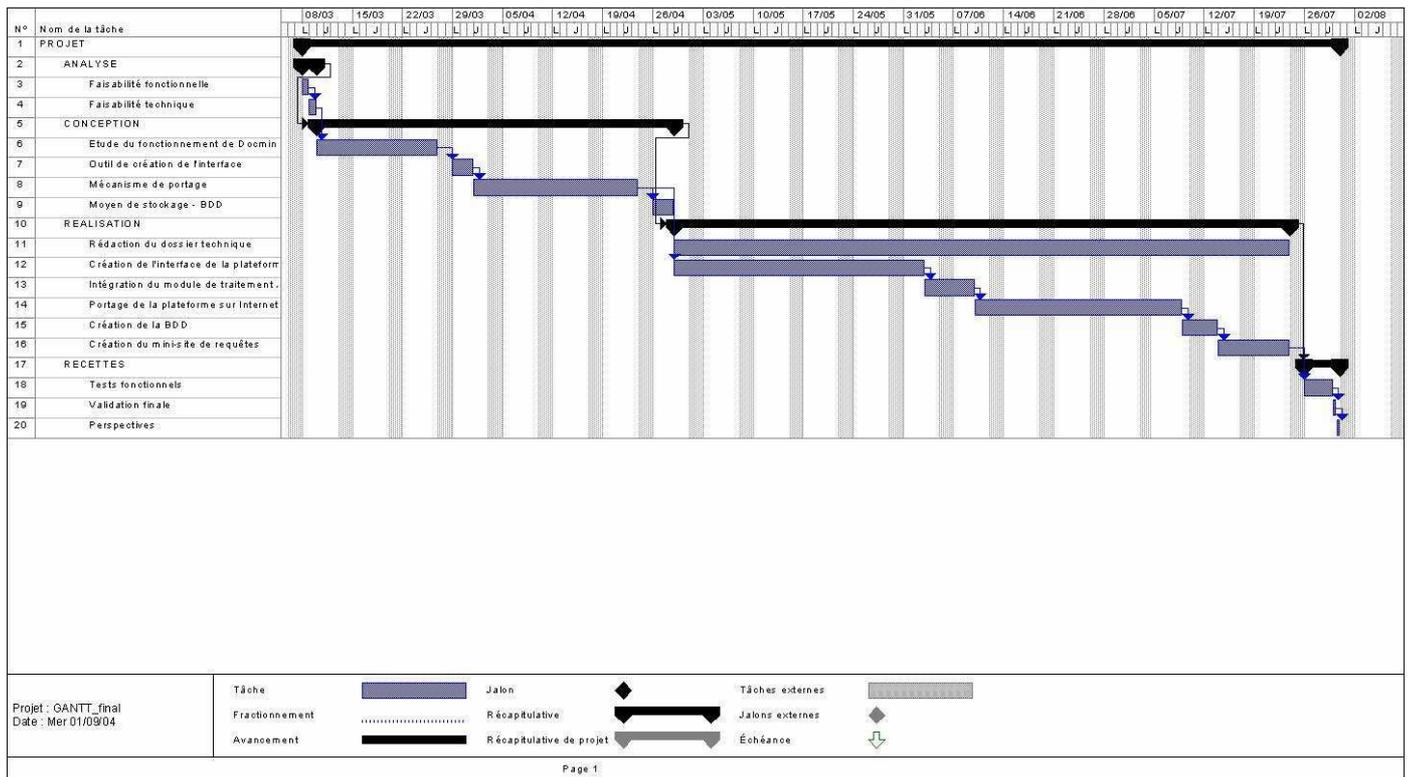


Figure 7 - Diagramme de GANTT

Durant ce travail de planification, plusieurs tâches principales ont été identifiées :

- analyse
- conception
- réalisation
- recette

Le diagramme de GANTT ci-dessus présente la subdivision du stage. Chaque tâche principale est décomposée en plusieurs sous tâches telles que *Faisabilité fonctionnelle et technique* pour la phase d'Analyse ou *le moyen de stockage* pour la phase de Conception. Concernant le suivi du projet, un seul jalon de validation a été fixé : la validation finale du stage.

Un suivi de planification et d'avancement du projet dans le temps nous a amené à réaliser un diagramme de GANTT supplémentaire (fourni en Annexes). Ce dernier montre l'état d'avancement de chaque tâche et de la totalité du stage. Cela nous permet de visualiser plus rapidement le possible retard ou avancement.

B- Risque du projet

Au cours de l'analyse du sujet du stage, trois risques ont identifiés :

- la non disponibilité ou l'absence de Mr OGIER au moment opportun
- une difficulté de compréhension du fonctionnement du projet Docmining (dû à l'obligation de recréer la partie Interface et de l'intégrer à la partie Traitement)
- la possible perte de temps à créer ce module Interface

Ils ont été suivis de manière rigoureuse, par le biais de trois fiches de risques, pour prévenir leur apparition. Ces dernières expliquent de quelle manière un suivi de ces risques doit être réalisé. Elles décrivent également les actions préventives et curatives à mettre en place en cas d'apparition et montrent l'évaluation de leur criticité.

Pour éviter leur apparition, toutes les actions décrites précédemment sont mises en œuvre.

<p>Stage : Projet M.A.D.O.N.N.E</p>	<p>Nom du risque : Difficulté de compréhension du fonctionnement du projet Docmining</p>	<p>Fiche de suivi de risque N° : 2 / version : 1</p>																																																																																															
<p>Date d'identification : 15 mars 2004 Date de mise à jour : 15 mars 2004 Responsable suivi : Wilfrid Papin</p>	<p>Description du risque : La création de la partie Interface de la plateforme SOX et de son intégration au module Traitement m'a obligé à travailler sur le projet initiateur du projet M.A.D.O.N.N.E, le projet Docmining. Une difficulté à comprendre son fonctionnement ferait perdre beaucoup de temps pour la suite du stage</p>																																																																																																
<p>Evaluation de la criticité</p> <table border="1" style="width: 100%; border-collapse: collapse; text-align: center;"> <tr> <td style="padding: 2px;">Gravité</td> <td style="padding: 2px;">0</td> <td style="padding: 2px;">1</td> <td style="padding: 2px;">2</td> <td style="padding: 2px;">3</td> </tr> <tr> <td style="padding: 2px;">Délai</td> <td></td> <td></td> <td></td> <td style="text-align: right;">X</td> </tr> <tr> <td style="padding: 2px;">Coûts</td> <td></td> <td></td> <td style="text-align: right;">X</td> <td></td> </tr> <tr> <td style="padding: 2px;">Recettes</td> <td></td> <td style="text-align: right;">X</td> <td></td> <td></td> </tr> <tr> <td style="padding: 2px;">Performances techniques</td> <td></td> <td></td> <td></td> <td style="text-align: right;">X</td> </tr> <tr> <td style="padding: 2px;">Autre :</td> <td></td> <td></td> <td></td> <td></td> </tr> <tr> <td style="padding: 2px;">Globale</td> <td></td> <td></td> <td></td> <td style="text-align: right;">X</td> </tr> </table> <div style="margin-top: 5px;"> <table border="1" style="width: 100%; border-collapse: collapse; text-align: center;"> <tr> <td style="padding: 2px;">G</td> <td style="padding: 2px;">3</td> <td style="padding: 2px;">2</td> <td style="padding: 2px;">1</td> <td style="padding: 2px;">0</td> </tr> <tr> <td style="padding: 2px;">R</td> <td style="padding: 2px;">3</td> <td style="padding: 2px;">2</td> <td style="padding: 2px;">1</td> <td style="padding: 2px;">0</td> </tr> <tr> <td style="padding: 2px;">A</td> <td style="padding: 2px;">3</td> <td style="padding: 2px;">2</td> <td style="padding: 2px;">1</td> <td style="padding: 2px;">0</td> </tr> <tr> <td style="padding: 2px;">V</td> <td style="padding: 2px;">3</td> <td style="padding: 2px;">2</td> <td style="padding: 2px;">1</td> <td style="padding: 2px;">0</td> </tr> <tr> <td style="padding: 2px;">I</td> <td style="padding: 2px;">3</td> <td style="padding: 2px;">2</td> <td style="padding: 2px;">1</td> <td style="padding: 2px;">0</td> </tr> <tr> <td style="padding: 2px;">T</td> <td style="padding: 2px;">3</td> <td style="padding: 2px;">2</td> <td style="padding: 2px;">1</td> <td style="padding: 2px;">0</td> </tr> <tr> <td style="padding: 2px;">É</td> <td style="padding: 2px;">3</td> <td style="padding: 2px;">2</td> <td style="padding: 2px;">1</td> <td style="padding: 2px;">0</td> </tr> </table> <p style="text-align: center; margin-top: 5px;">PROBABILITE</p> <table border="1" style="width: 100%; border-collapse: collapse; text-align: center;"> <tr> <td style="padding: 2px;">Probabilité d'occurrence</td> <td style="padding: 2px;">3</td> </tr> <tr> <td style="padding: 2px;">0 : nulle</td> <td style="padding: 2px;">1 : <5%</td> </tr> <tr> <td style="padding: 2px;">2 : 5% < <20%</td> <td style="padding: 2px;">3 : >20%</td> </tr> </table> </div>	Gravité	0	1	2	3	Délai				X	Coûts			X		Recettes		X			Performances techniques				X	Autre :					Globale				X	G	3	2	1	0	R	3	2	1	0	A	3	2	1	0	V	3	2	1	0	I	3	2	1	0	T	3	2	1	0	É	3	2	1	0	Probabilité d'occurrence	3	0 : nulle	1 : <5%	2 : 5% < <20%	3 : >20%	<table border="1" style="width: 100%; border-collapse: collapse;"> <thead> <tr> <th style="padding: 2px;">Actions (préventif et/ou curatif)</th> <th style="padding: 2px;">Qui</th> <th style="padding: 2px;">Date début</th> <th style="padding: 2px;">Date fin</th> </tr> </thead> <tbody> <tr> <td style="padding: 2px;"><u>Actions préventives</u></td> <td></td> <td></td> <td></td> </tr> <tr> <td style="padding: 2px;">- Pour perdre le moins de temps possible, contacter le plus rapidement possible les concepteurs (dès qu'il y a une incompréhension)</td> <td></td> <td style="text-align: center;">15-03-2004</td> <td style="text-align: center;">30-07-2004</td> </tr> <tr> <td style="padding: 2px;"><u>Actions curatives</u></td> <td></td> <td></td> <td></td> </tr> <tr> <td style="padding: 2px;">- Si perte de temps importantes, discussion avec Jean Marc Ogier pour trouver la marche à suivre</td> <td></td> <td></td> <td></td> </tr> </tbody> </table>	Actions (préventif et/ou curatif)	Qui	Date début	Date fin	<u>Actions préventives</u>				- Pour perdre le moins de temps possible, contacter le plus rapidement possible les concepteurs (dès qu'il y a une incompréhension)		15-03-2004	30-07-2004	<u>Actions curatives</u>				- Si perte de temps importantes, discussion avec Jean Marc Ogier pour trouver la marche à suivre			
Gravité	0	1	2	3																																																																																													
Délai				X																																																																																													
Coûts			X																																																																																														
Recettes		X																																																																																															
Performances techniques				X																																																																																													
Autre :																																																																																																	
Globale				X																																																																																													
G	3	2	1	0																																																																																													
R	3	2	1	0																																																																																													
A	3	2	1	0																																																																																													
V	3	2	1	0																																																																																													
I	3	2	1	0																																																																																													
T	3	2	1	0																																																																																													
É	3	2	1	0																																																																																													
Probabilité d'occurrence	3																																																																																																
0 : nulle	1 : <5%																																																																																																
2 : 5% < <20%	3 : >20%																																																																																																
Actions (préventif et/ou curatif)	Qui	Date début	Date fin																																																																																														
<u>Actions préventives</u>																																																																																																	
- Pour perdre le moins de temps possible, contacter le plus rapidement possible les concepteurs (dès qu'il y a une incompréhension)		15-03-2004	30-07-2004																																																																																														
<u>Actions curatives</u>																																																																																																	
- Si perte de temps importantes, discussion avec Jean Marc Ogier pour trouver la marche à suivre																																																																																																	
<p>Suivi proposé</p> <p>Rester en contact permanent avec les concepteurs du projet Docmining ainsi que ses utilisateurs afin de ne pas perdre trop de temps</p>	<p>Commentaires / Bilan</p>																																																																																																
<table style="width: 100%;"> <tr> <td style="width: 40%; vertical-align: top;"> <p>Gravité du risque</p> <p>0 - aucun impact sur les objectifs du projet</p> <p>1 - objectifs du projet affectés « à la marge » (domaine du Chef de Projet)</p> <p>2 - rentabilité du projet significativement affectée</p> <p>3 - le projet est condamné</p> </td> <td style="width: 20%; text-align: center; vertical-align: middle;"> <p>Légende</p> </td> <td style="width: 40%; vertical-align: top;"> <p>Criticité du risque</p> <p>Gravité ou probabilité d'occurrence nulle ou mineure</p> <p>Problème technique courant</p> <p>Remet en cause les solutions fondamentales</p> <p>Inacceptable : incompatible avec le besoin exprimé</p> </td> </tr> </table>			<p>Gravité du risque</p> <p>0 - aucun impact sur les objectifs du projet</p> <p>1 - objectifs du projet affectés « à la marge » (domaine du Chef de Projet)</p> <p>2 - rentabilité du projet significativement affectée</p> <p>3 - le projet est condamné</p>	<p>Légende</p>	<p>Criticité du risque</p> <p>Gravité ou probabilité d'occurrence nulle ou mineure</p> <p>Problème technique courant</p> <p>Remet en cause les solutions fondamentales</p> <p>Inacceptable : incompatible avec le besoin exprimé</p>																																																																																												
<p>Gravité du risque</p> <p>0 - aucun impact sur les objectifs du projet</p> <p>1 - objectifs du projet affectés « à la marge » (domaine du Chef de Projet)</p> <p>2 - rentabilité du projet significativement affectée</p> <p>3 - le projet est condamné</p>	<p>Légende</p>	<p>Criticité du risque</p> <p>Gravité ou probabilité d'occurrence nulle ou mineure</p> <p>Problème technique courant</p> <p>Remet en cause les solutions fondamentales</p> <p>Inacceptable : incompatible avec le besoin exprimé</p>																																																																																															

Figure 8 - Feuille de risque

Tout au long de ce stage, les deux risques identifiés sont suivis par le biais de feuille de suivi de risques permettant de voir l'évolution du risque, sa criticité et permettent de limiter leur apparition en les surveillant de manière continue.

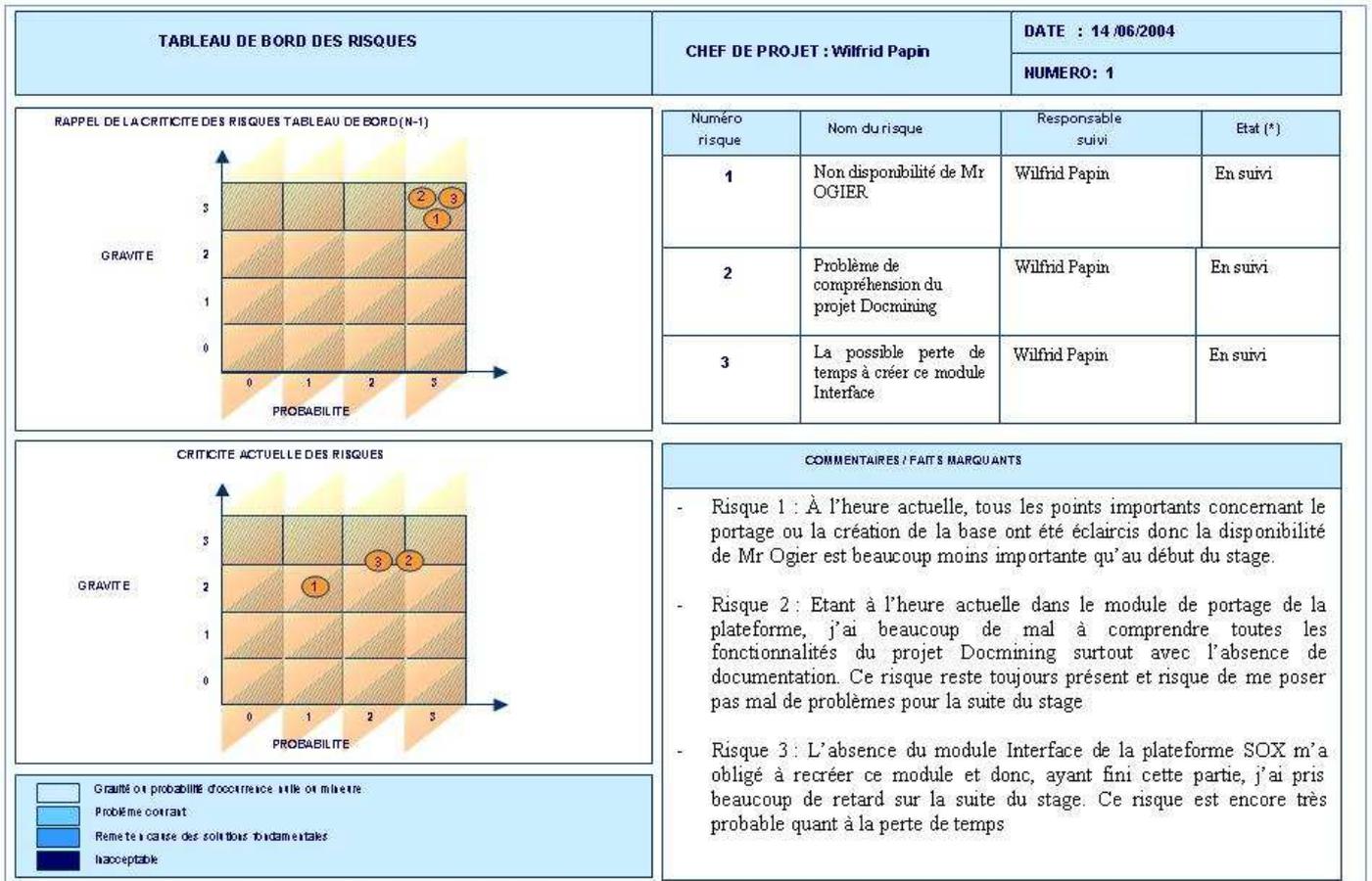


Figure 9 - Feuille de suivi des risques

C- Consortium et réunions

L'avancement du projet est ponctué de réunions internes organisées entre Jean Marc OGIER et moi-même pour trouver des solutions adaptées aux problèmes mais également pour analyser les points à améliorer afin d'obtenir un bon résultat.

Un consortium a également été organisé le 25 Mai à Tours avec tous les partenaires du projet M.A.D.O.N.N.E afin de faire un état de l'art du travail de chacun. Une présentation du stage a permis de voir les points qui n'allaient pas et ceux qui répondaient parfaitement à la demande de chacun.

VI- Environnement technique

A- Architecture logicielle

L'infrastructure informatique mise en place au laboratoire a deux objectifs :

- la mise à disposition d'outils informatiques à usage interne pour le personnel ;
- l'ouverture de ses ressources vers l'extérieur (monde recherche) afin de promouvoir certains travaux comme par exemple le projet MADONNE;

Le schéma ci-dessous identifie à la fois les moyens matériels et logiciels nécessaires au fonctionnement des différents outils utilisés dans le laboratoire. Cette vue de l'architecture logicielle est complètement intégrée au reste de l'infrastructure informatique de l'Université de La Rochelle.

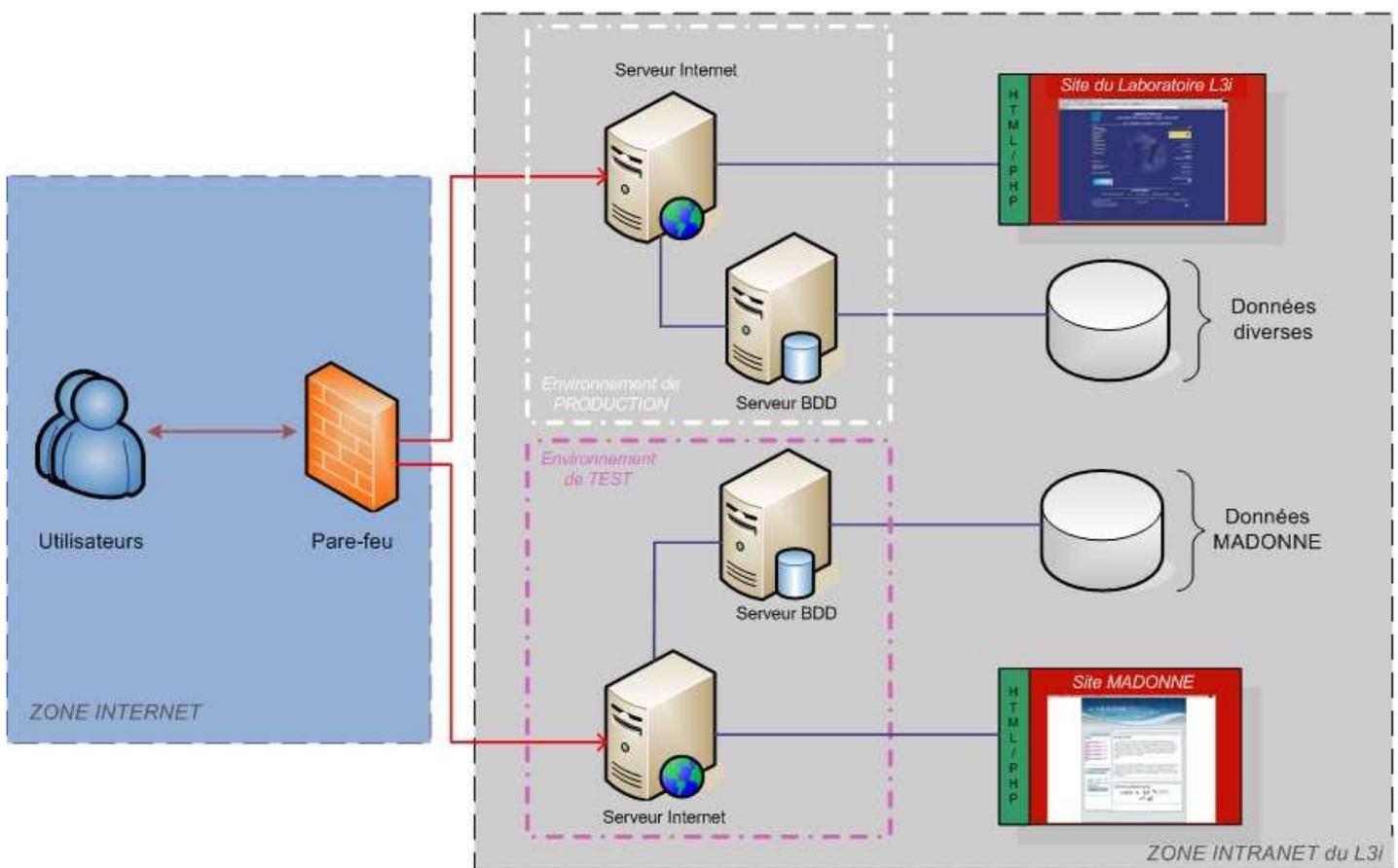


Figure 10 - Infrastructure matérielle et logicielle

Deux mondes distincts sont identifiables : la zone internet et la zone intranet du laboratoire L3i. Le lien entre ces deux espaces est réalisé par un pare-feu destiné à assurer la protection des moyens informatiques internes. Il existe d'autres serveurs internes au laboratoire mais ils ne servent que de serveurs de test pour certaines applications et même de serveur de stockage uniquement.

B- Données existantes

En ce qui concerne la partie Base de données, seul les images existaient. Elles étaient présentes sur certains ordinateurs et sur le serveur FTP du Centre d'Études Supérieures de la Renaissance. Donc à ce niveau là, un travail important de récupération est à faire en plus de la création de la base.

Pour ce qui est de la partie Portage de la plateforme SOX sur Internet, certains faits amènent à réévaluer le travail à effectuer.

En effet, lors de mon arrivée, la plateforme SOX était sans savoir pourquoi introuvable. Donc après quelques recherches et une discussion entre Jean Marc OGIER et moi-même, la décision de refaire entièrement la partie Interface de la plateforme a été prise. Heureusement, la partie Traitement de cette plateforme n'était pas à refaire car le projet SOX n'était que la suite logique du projet Docmining*. Finalement, en plus du portage sur le Web, la partie interface est à créer avec intégration de la partie traitement créée il y a déjà quelques temps.

(Docmining* : projet exploratoire dont le but est de définir des services d'intermédiation dans l'échange de documents hétérogènes incluant des composantes non textuelles au sens ASCII du terme.)

VII- Conception

A- Architecture du stage

Les fonctionnalités demandées dans le sujet de stage nécessitent la conception de trois *briques* logicielles complémentaires. Le schéma ci-dessous illustre leurs interactions :

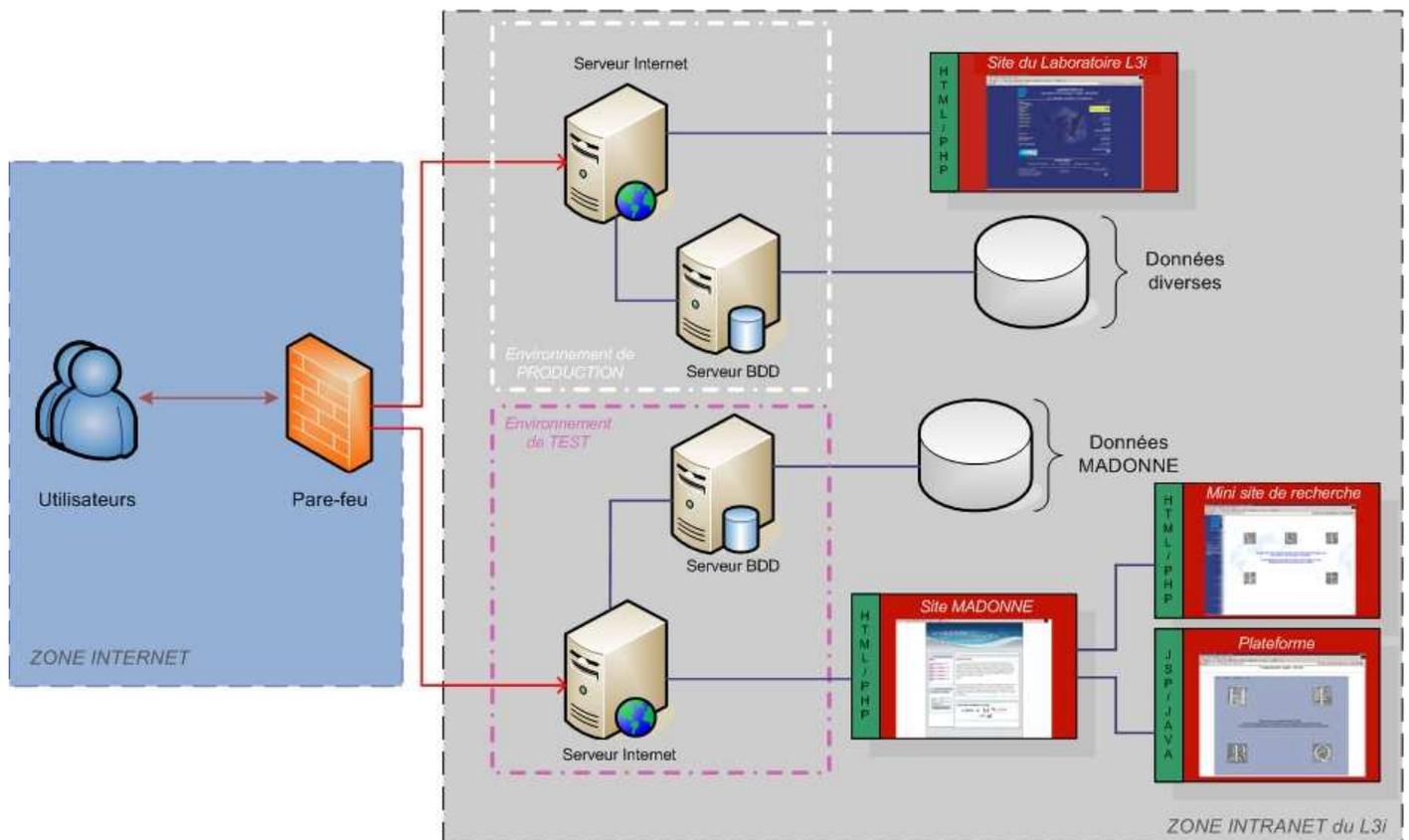


Figure 11 - Nouvelle infrastructure avec prise en compte du travail réalisé

La première a pour rôle de créer un mini site permettant à un utilisateur de faire une recherche d'image sur la base de données. La nature de la recherche est très différente de celle que l'on connaît puisque la navigation ne se fait pas forcément par mot-clef. Elle peut se faire dans notre cas soit par méta-données, soit par un calcul de distance entre une signature de référence et celles des images de la base (voir les détails dans la partie Réalisation, page 49).

Pour ce faire, il est nécessaire de créer des pages en HTML et PHP interrogeant la base via des requêtes créées en SQL. Cette extension est identifiée sur le schéma ci-dessus par *Mini site de recherche*.

La seconde brique identifiée est l'élément de fondation pour le mini site. Un espace de stockage pour les images permet à la fois de recevoir les informations provenant des requêtes mais permet également de répondre aux requêtes des utilisateurs. Du point de vue informatique, cela est réalisé par la création d'une base de données et l'ajout de différentes tables permettant la gestion des images ainsi que les méta-données et les signatures associées. Sur le schéma cela est intitulé *Données Images*.

La dernière brique constitue l'autre partie du stage, le développement de la plateforme logicielle ainsi que son portage sur Internet. Ce module est réalisé en Java, Java Swing et surtout en Servlet pour le portage. Cette technologie nécessite un serveur Internet pour son bon fonctionnement donc cette brique est directement intégrée au site MADONNE. Sur le schéma, il est représenté par *Plateforme MADONNE*.

Les parties qui suivent détaillent le fonctionnement de la plateforme SOX ainsi que la conception des *briques* logicielles décrites au-dessus.

B- Etude et fonctionnement de Docmining / SOX

1- Plateforme SOX : Fonctionnement

N'ayant pu ni voir, ni étudier cette plateforme SOX, je vais vous en décrire les objectifs ainsi que son utilisation de manière assez succincte.

L'objectif de SOX est de constituer une plate-forme pour le traitement de documents. Cette plate-forme comporte différentes briques logicielles :

- Des bibliothèques de traitements provenant des différents membres du consortium ou de logiciels libres.
- Un moteur d'interprétation de scénario de traitements.
- Des outils pour la définition de scénario.
- Des outils pour intégrer de nouveaux traitements à partir de logiciels externes.
- Des outils pour intégrer de nouveaux traitements à partir de scénarios SOX.
- Des outils pour la visualisation des résultats d'un scénario.
- Divers outils de manipulation de fichiers image et XML.

Un traitement dans SOX est défini par son fichier XML de propriété indiquant :

- Dans quelle classe Java il est implémenté.
- Le(s) service(s) fourni(s)
- Les objets manipulés.
- Les objets produits à partir des objets manipulés.

Cette plateforme intègre des traitements de langage différents comme Java mais aussi C, C++. La communication entre des traitements en C/C++ et le langage utilisé pour développer l'interface, Java, se fait par les JNI (Java Native Interface).

Son fonctionnement est principalement basé sur les scénarii. L'utilisateur a la possibilité de créer, modifier des scénarii afin d'exécuter des traitements hétérogènes comme la binarisation, la décomposition en composantes connexes, etc...

Ces traitements sont applicables sur des documents ayant des formats très différents (Pdf, PostScript, Pgm, etc...).

Concernant la partie Interface de cette plateforme, elle donne la possibilité à l'utilisateur de créer, modifier, dupliquer, exécuter des scénarii. Ensuite, le résultat de l'exécution d'un traitement sur un document se visualise dans une «double» interface nommée Xmillum et XML Tree. La première interface permet de visualiser graphiquement le résultat et la deuxième permet de le visualiser via un arbre XML.

La partie suivante décrit le problème survenu au niveau de la plateforme SOX.

2- Problème majeur sur la plateforme SOX

Lors de mon arrivée au laboratoire, mon travail a commencé par une étude du projet Docmining pour bien comprendre le fonctionnement des traitements. Par la suite, je devais faire l'étude de la plateforme SOX afin de trouver les meilleurs outils pour son portage sur Internet.

Mais mon étude de la plateforme SOX n'a pu débuter car mon maître de stage était dans l'impossibilité de la trouver. Après une recherche infructueuse sur les serveurs du laboratoire, dans les archives de Mr OGIER et auprès de l'étudiant ayant réalisé la plateforme, la décision de refaire la partie Interface de la plateforme a été prise. Les parties Traitement et Visualisation de la plateforme SOX étant quasiment identique à celle de la plateforme Docmining, ces deux parties ne sont pas à refaire car la plateforme Docmining étant présente sur les serveurs du laboratoire.

3- Plateforme Docmining : Fonctionnement

Le fonctionnement de la plateforme Docmining est équivalent à celui de SOX. La différence se fait essentiellement au niveau de l'interface car le projet Docmining ne possède pas d'interface de présentation et d'exécution de traitements. Il y a seulement l'interface de visualisation (Xmillum et SCI Scénario Construction Interface). Il y a quand même une légère différence au niveau de l'interface de visualisation, il s'agit de l'arbre XML. Sur la plateforme SOX, il s'agit d'un arbre XML normal regroupant tous les objets rencontrés dans le document traité. Par contre, sur la plateforme Docmining, il s'agit seulement d'un arbre d'objets mais ayant la même utilisation et le même rendu.

Sinon, pour lancer l'exécution d'un traitement sur un document, il suffit de le faire en ligne de commande :

```
C:\java ScenarioConstruction intrac.xml figaro.xml visucc.xsl
```

Le fichier ScenarioConstruction est le fichier contenant la procédure principale. Le fichier figaro.xml est le fichier contenant le chemin vers le document à traiter. Intrac.xml est le fichier de paramètres indiquant le chemin vers les fichiers XML de traitements et le fichier visucc.xsl indiquant quels types de fichiers sont pris en compte dans l'exécution.

Pour le résultat, il est équivalent à la plateforme SOX car le document est visible sur la partie de l'interface Xmillum. Ensuite, l'utilisateur sélectionne le document dans l'arbre d'objet comme sur l'image suivante.

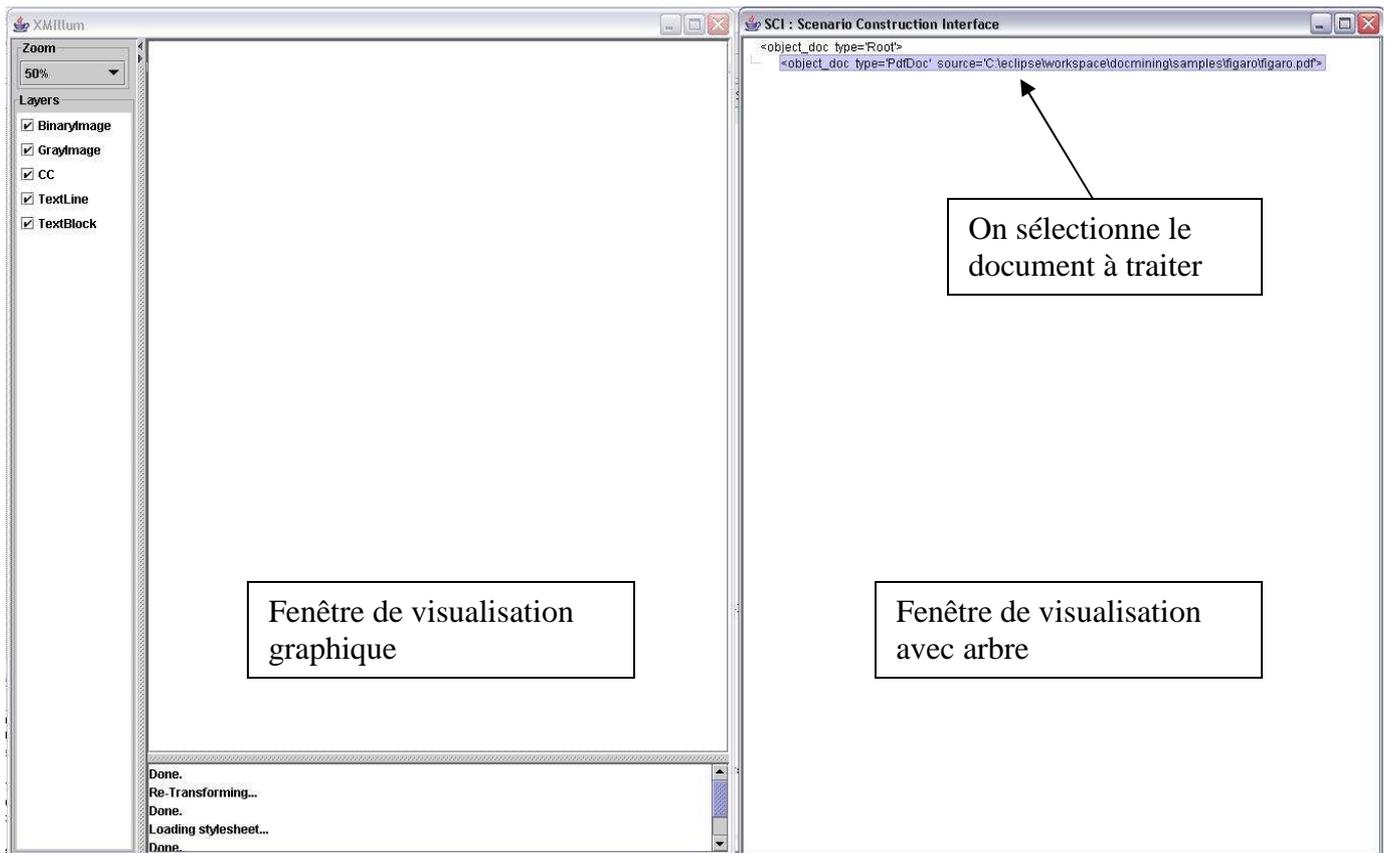


Figure 12 - Fenêtre de visualisation XMillum

Après avoir sélectionner le document, l'utilisateur a la possibilité, via le clic droit de la souris, d'exécuter un traitement et de voir le résultat dans la fenêtre Xmillum.

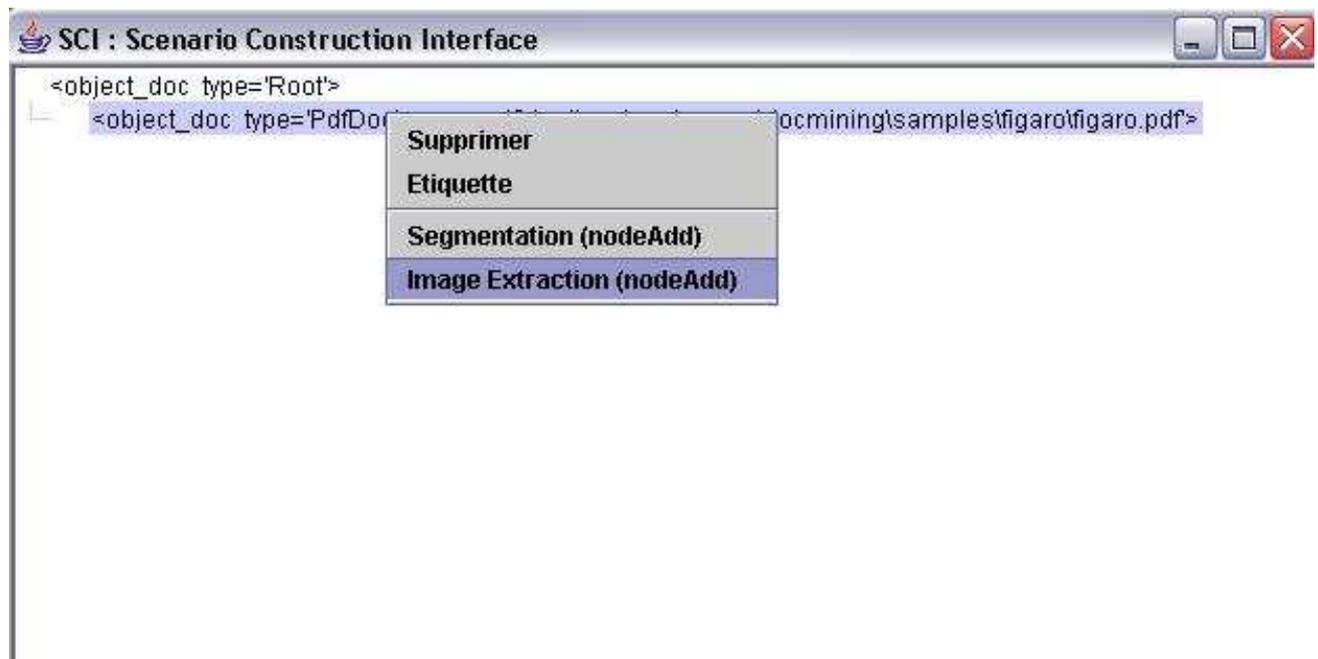


Figure 13 - Sélection du traitement à exécuter



Figure 14 - Résultat du traitement

L'exécution d'un traitement passe d'abord par la validation d'une fenêtre de paramètre.

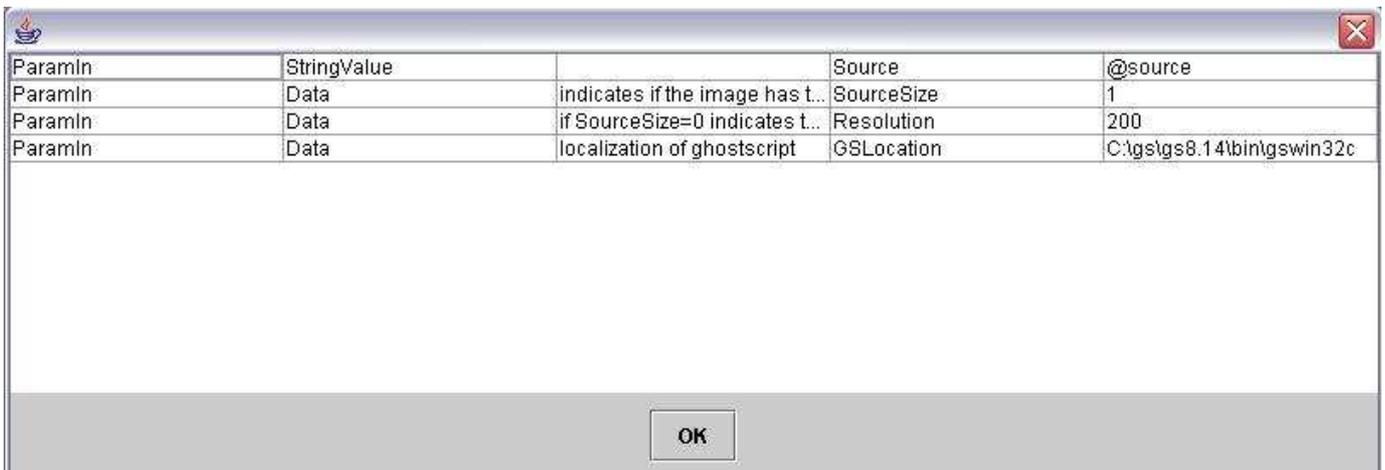


Figure 15 - Fenêtre de paramètre

La partie suivante décrit les outils choisis pour la création de l'interface.

C- Outil de création pour l'interface

Le premier travail pour la création de l'interface consiste à identifier toutes les fonctionnalités nécessaires à une utilisation adéquate. Comme cette interface doit reprendre celle de la plateforme SOX, les fonctionnalités doivent être identiques.

Après l'étude du fonctionnement de la plateforme Docmining et comme l'ensemble des fichiers de visualisation sont en java, l'outil le plus adéquate pour la création de la plateforme vu les techniques apprises au cours de mon cursus, est le Java Swing.

Le Java Swing est une bibliothèque de classes permettant de créer des interfaces Utilisateur puis de les gérer. Il existe une seconde bibliothèque de classes permettant la même chose, l'AWT. La procédure de construction d'une interface Swing est similaire à celle d'AWT : créer un cadre, des composants dans ce cadre, une mise en page pour ces composants, des méthodes en réponses aux actions de l'utilisateur.

Mais la différence entre Swing et AWT est majeure du côté Utilisateur : l'apparence des composants est totalement différente. Du côté Concepteur, Swing présente plus de composants qu'AWT donc plus de possibilité de créer des interfaces complexes et efficaces. La gestion des événements est également complètement différente.

Toutes ces différences ont fait de la bibliothèque Swing le choix pour la création de l'interface de la plateforme SOX.

Suite à l'étude du fonctionnement de la plateforme SOX et de celle de Docmining, certains objets de la bibliothèque Swing ont été choisis pour faire de l'interface SOX une interface facile d'utilisation mais très complète pour des utilisateurs spécialistes.

Les objets choisis sont par exemple :

- JTree : arbre dynamique
- JList : liste d'objets
- JTable : tableau
- Etc...

Pour conclure, la plateforme, dont les outils viennent d'être décrits précédemment, doit être utilisable sur Internet. Il y a donc plusieurs possibilités pour ce portage. Mais concernant l'interface, le choix est simple. Étant donné l'utilisation de ce module, le choix de créer l'interface comme un applet a été pris. Un applet étant une application intégrée dans une page internet. Il s'agit d'une portion de code Java résidant sur un serveur et apte à être téléchargée sur un navigateur Web doté d'une machine virtuelle Java (interpréteur de code Java). L'avantage d'un applet est d'offrir à l'utilisateur du navigateur une interface beaucoup plus interactive du même niveau que l'interface graphique du système d'exploitation qu'une page HTML munie de ses seuls liens hypertextes.

La partie suivante décrit la conception du mécanisme de portage de la plateforme sur Internet.

D- Mécanisme de portage

Ce mécanisme de portage est en fait une partie déjà dans le module de création de l'interface. En effet, le choix de la créer par le biais d'une applet montre déjà ce mécanisme de portage. Le fait de considérer une interface comme une applet n'est due qu'à une extension, « JApplet ». La classe principale permettant de construire cette interface utilise des outils de base comme les JButton, JLabel, etc... Pour porter la plateforme sur Internet, la classe principale doit instancier la classe JApplet transformant donc l'interface en applet.

Mais il y a un problème important dû à l'utilisation d'une applet, les restrictions de sécurité. En effet, une applet a des restrictions :

- Elle ne peut lire un fichier présent sur le disque dur de celui la lançant
- Elle ne peut non plus écrire sur le disque dur de celui la lançant
- Elle ne pas se connecter via des sockets à un host différent de celui l'ayant téléchargée

Et donc en signant l'applet, ces restrictions sont oubliées. Afin de signer l'applet, il faut d'abord créer un certificat. Pour se faire, il faut générer une clé publique via l'outil présent dans le JDK, « keytool.exe ». Toutes ces clés seront stockées dans la « keystore ». La ligne de commande permettant cela est :

keytool -genkey -alias cle_app

- genkey : active la génération des clés publique et privée
- alias : permet de donner un nom à la clé

Lors de la création de cette clé, plusieurs paramètre sont à renseignés comme un mot de passe pour la keystore, un nom et prénom du signataire, nom de l'organisation, de la ville, du pays et le code du pays.

Après la création du certificat, l'archive doit également être signer en utilisant ce dernier via la ligne de commande suivante :

jarsigner -verbose nom_archive.jar myApp

Afin de vérifier la mise sur le fichier de la signature, la commande suivante est à utiliser :

jarsigner -verify -verbose -certs nom_archive.jar

Pour conclure cette procédure de création d'un certificat, il faut générer ce dernier. Pour cela, la commande est la suivante :

keytool -export -alias cle_app -file myApp.cer

- export : permet de lire la clé associé à cle_app
- file : nom du fichier de sortie contenant le certificat (extension .cer)

Cette applet s'utilise donc dans un navigateur Internet et correspond à la partie visible de la plateforme. Pour l'utiliser sur un navigateur, l'archive de l'applet doit se trouver sur le serveur Web. Par la suite, une requête est proposée au serveur par l'applet suite à l'utilisation de l'applet par un utilisateur. Cette communication entre le serveur et l'applet doit se gérer par une application sachant comprendre les requêtes HTTP et pouvant envoyer une réponse HTTP.

Après une recherche des outils les plus adéquates et les plus efficaces afin de gérer cette communication, deux techniques ont été sélectionnées :

- les servlets
- le framework Struts

Les servlets sont des applications Java fonctionnant du côté serveur. Elles permettent donc de gérer des requêtes HTTP et de fournir au client une réponse HTTP dynamique. Etant donné qu'il s'agit d'une technologie Java, les servlets fournissent un moyen d'améliorer les serveurs web sur n'importe quelle plateforme, d'autant plus que les servlets sont indépendants du serveur Web.

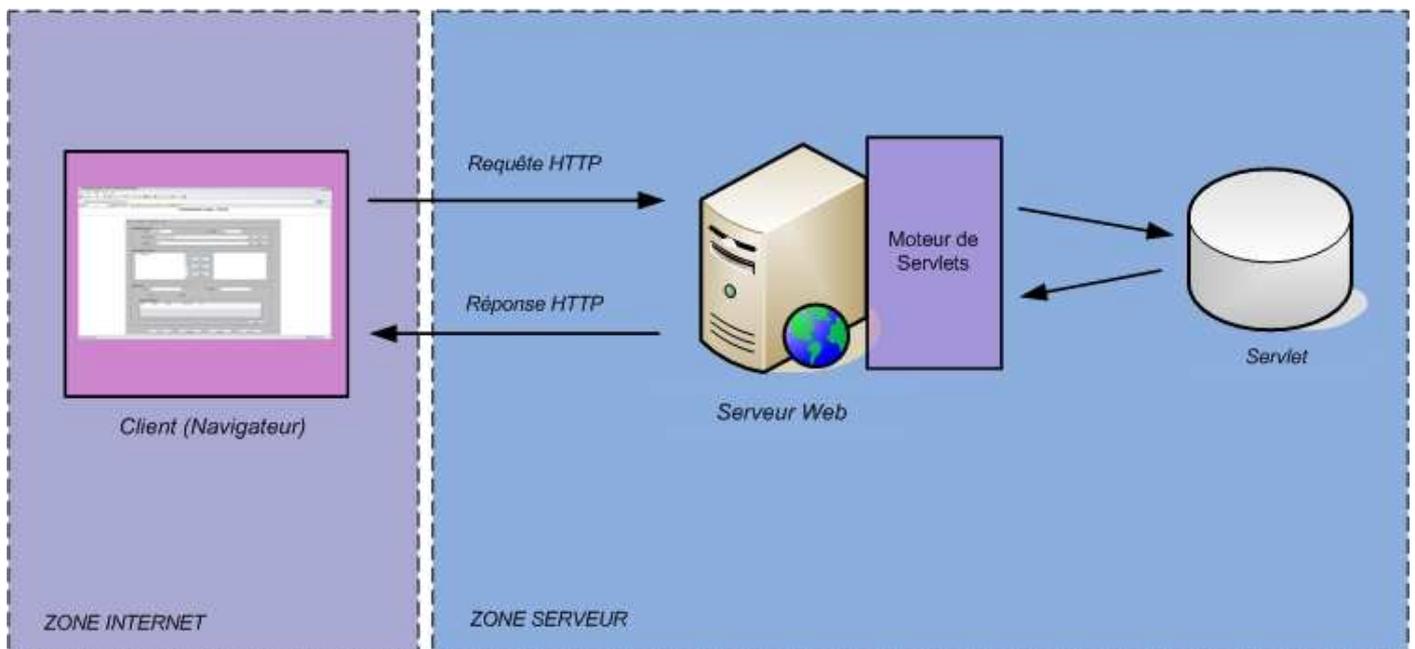


Figure 16 - Fonctionnement d'une servlet

Concernant Struts, ce framework propose un cadre logiciel pour l'organisation des échanges avec l'utilisateur et la dynamique de l'application.

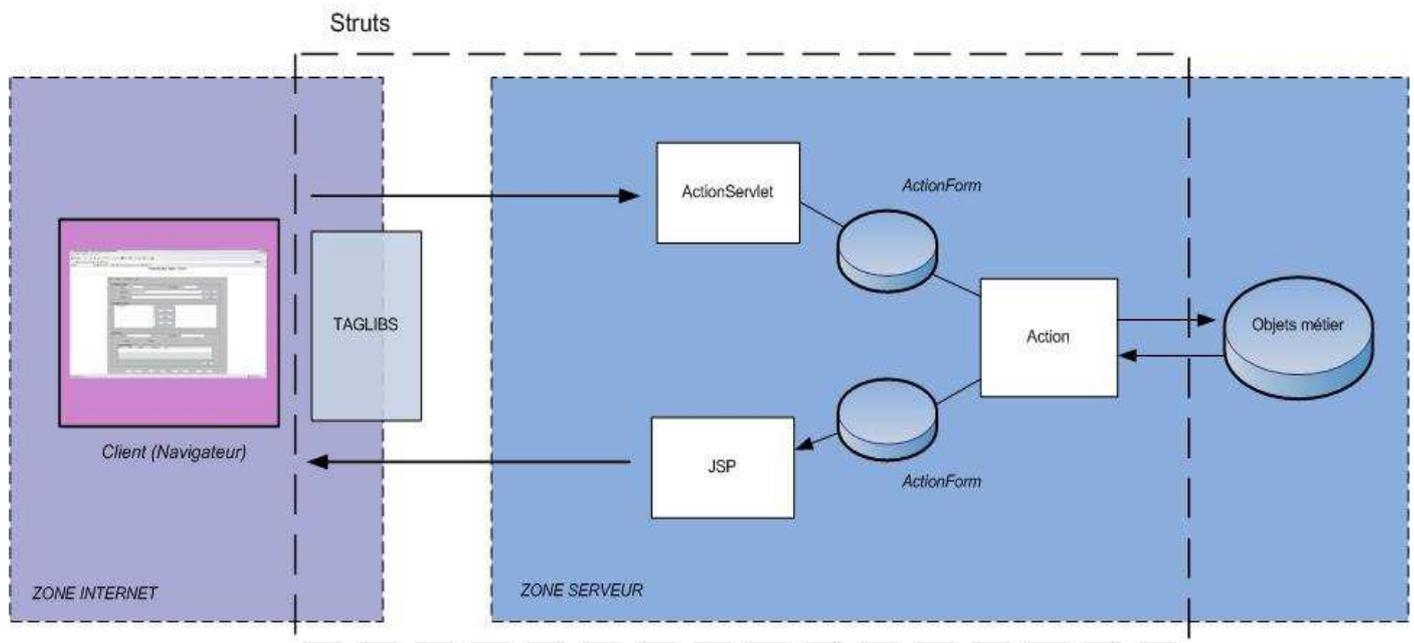


Figure 17 - Fonctionnement de Struts

Struts s'appuie sur le modèle de conception des commandes : un discriminant, une chaîne de caractères en général, fait office de commande et indique le traitement à exécuter. Dans Struts, la commande est déterminée par une partie de l'URL transmise au servlet contrôleur, *ActionServlet*. Les traitements sont encapsulés dans les classes *Actions*. La correspondance entre la commande et l'*Action* à exécuter est définie par les classes *ActionMapping*.

Struts intègre un mécanisme d'automatisation des relations entre *Actions* et *JSP* par le biais des *ActionForms*. Ces classes sont des *JavaBeans* renseignés automatiquement à partir des paramètres des requêtes transmises au servlet contrôleur et transmises telles quelles aux *Actions*.

Après avoir fait une étude assez précise par le biais d'exemples, de tutoriaux, le choix de prendre la version la plus simple de la communication entre Applet et Serveur, les servlets, fût prise. Ce choix s'est fait également par rapport au travail à effectuer dans ce stage car le re-développement de l'interface de la plateforme a été un travail important et supplémentaire pesant sur ce choix.

La dernière partie expose la technique employée pour le moyen de stockage des images.

E- Moyen de stockage des images

Il existe deux méthodes permettant de stocker des images dans une base de données. La première est de stocker directement les images dans la base via des objets tels que les BLOB (Binary Large Object). Ce sont des objets dédiés aux données binaires et donc notamment aux images.

La seconde méthode est de mettre les images dans de simples fichiers organisés en répertoire et de gérer les chemins à ces fichiers dans la base de données.

Au final, après une étude assez rapide, la méthode avec les objets BLOB a été abandonnée car la manipulation de ces derniers ne peuvent pas la plupart du temps faire l'objet de requêtes SQL basiques. De plus, il ne faut pas oublier qu'un ordinateur est plus rapide pour la manipulation de fichiers et donc il y a plus à gagner à utiliser un stockage sous forme de fichiers que de manipuler des flux binaires dans des colonnes de type BLOB.

En revanche, les images possèdent des caractéristiques nommées « attributs » ou « méta-données » qu'il est souvent utile de collecter dans une base de données et la technique choisie ne gère pas ces attributs de manière automatique.

Les attributs les plus souvent utilisés sont :

FORMAT	BMP, JPEG, TIF
COULEURS	2 (noir & blanc), 256 (niveaux de gris) ...
NOM	Le nom du fichier
TAILLE	Le volume des données
DESCRIPTION	Une description de l'image pouvant faire l'objet d'une requête
COPYRIGHT	Auteur / propriétaire
DATE_CREATION	Date de création de l'image

En ce qui concerne le stockage des images, il est possible de mettre toutes les images dans un seul répertoire d'un serveur. Il suffit ensuite de connaître le chemin de ce répertoire que l'on nomme par exemple SIM_PATH. Dès lors, pour retrouver l'emplacement d'une image précise, la constante SIM_PATH est concaténée au nom de l'image.

Mais cette technique possède un inconvénient très important qu'il ne faut surtout pas oublier, il s'agit du trop grand nombre de fichiers dans un répertoire. En effet, les systèmes d'exploitations possèdent des limites et même si elles sont très larges, il n'est jamais impossible de les atteindre.

Il suffit donc de multiplier les répertoires et de stocker les images soit par type, par volume, par date de création. Et donc, pour savoir où chaque image se trouve, il suffit de rajouter aux caractéristiques déjà spécifiées, le chemin de stockage de chaque image. Ce chemin étant stocké dans une autre table.

Pour rajouter un plus au modèle expliqué au-dessus, ce dernier gère l'obsolescence des images. Autrement dit gérer une image qui n'est jamais utilisée ou encore remplacée par une autre. En ce qui concerne la suppression, une colonne supplémentaire est créée dans la table principale T_IMAGE_IMG afin d'y mettre la date de suppression. Et donc, toute requête pour retrouver une image pourra filtrer par cette date et celle du système.

Pour ce qui est du remplacement, il suffit de relier la table T_IMAGE_IMG à elle-même afin d'informer que l'image visée a été remplacée par une nouvelle. Une simple auto-référence suffit donc la colonne IMG_ID_REEMPLACEMENT a été créée pour cela.

A des fins de contrôles et de validation des données, une table de référence pour le format des images (T_TYPE_IMAGE_TIM) a été rajoutée au modèle. Cela permet de créer préventivement tous les types de formats avec lesquels les laboratoires partenaires désirent travailler.

Pour finir sur le moyen de stockage et afin de répondre totalement aux besoins de mon maître de stage, Mr OGIER et de ses partenaires, d'autres tables ont été créées. Ces besoins sont au niveau de la recherche multicritère (via la signature ou les méta-données). Donc 3 tables principales ont été rajoutées ainsi que 3 tables de « jointures » :

- la table T_OBJET : table correspondant aux méta-données qu'il est possible de rencontrer sur les images ;
- la table T_IMAGE_OBJET : table de jointure entre T_IMAGE_IMG et T_OBJET (dans le sens où une image peut contenir plusieurs objets et que un objet peut être contenu dans plusieurs images à la fois) ;
- la table T_SIGNATURE : table correspondant à une suite de valeurs que peut avoir une image ;
- la table T_IMAGE_SIGNATURE : table de jointure entre T_IMAGE_IMG et T_SIGNATURE (dans le sens où une image peut avoir plusieurs signatures)
- la table T_VALEUR : table correspondant à toutes les valeurs que peut avoir une signature
- la table T_VALEUR_SIGNATURE : table de jointure entre T_VALEUR et T_SIGNATURE (dans le sens où une signature peut avoir plusieurs valeurs différentes ou égales)

Voici donc le schéma final de cette base de données créée :

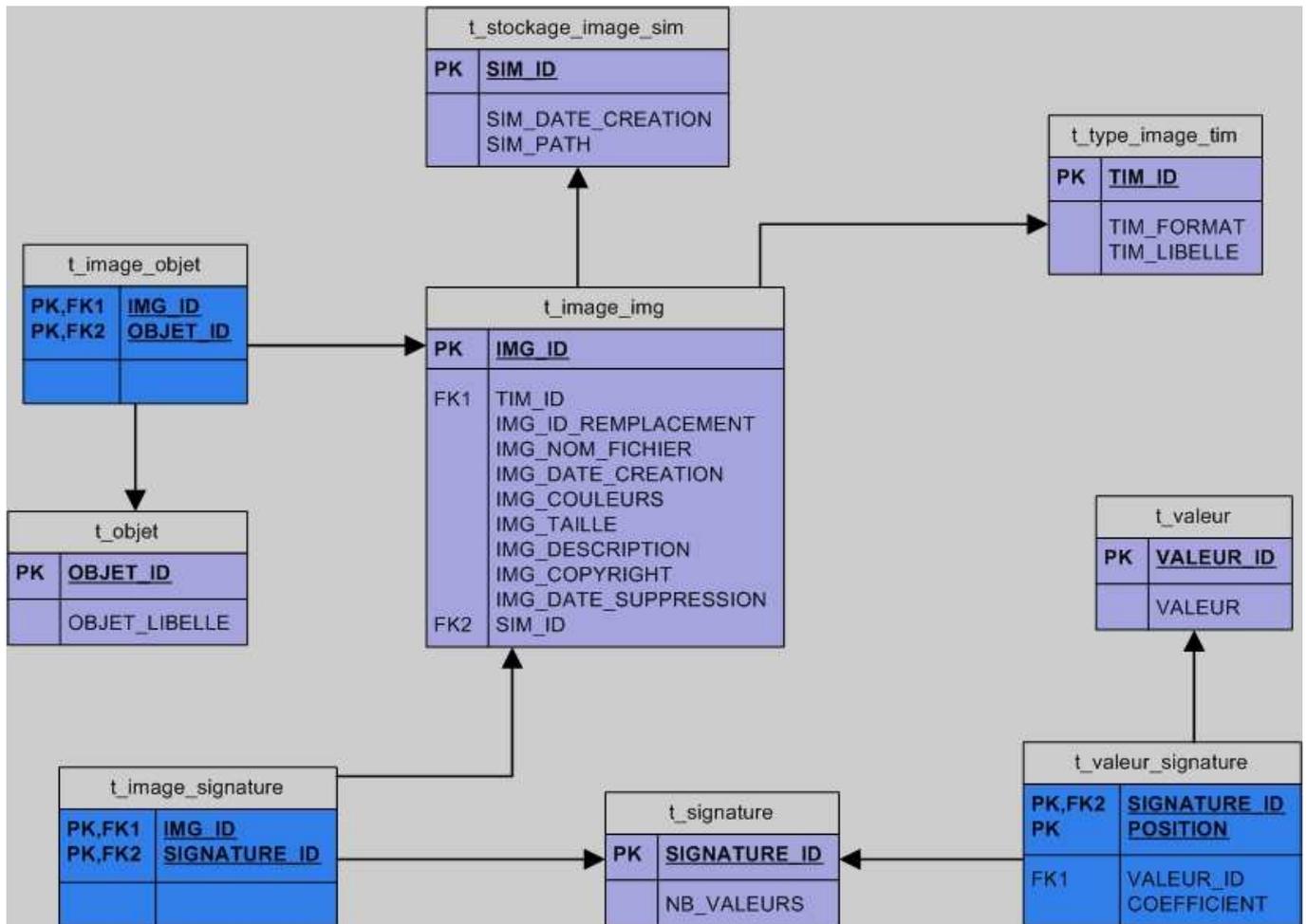
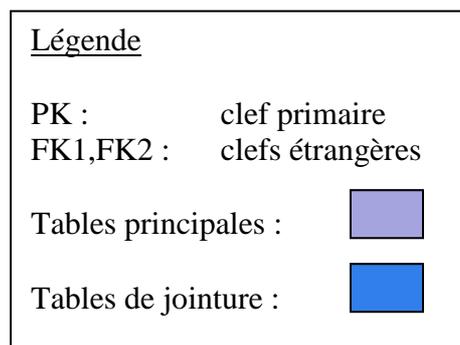


Figure 18 - Schéma de la base de données



Comme vous pouvez le voir sur le schéma ci-dessus, la table T_SIGNATURE a un champ NB_VALEURS correspondant au nombre de valeur qu'elle contient. Et pour ce qui est de la table T_VALEUR_SIGNATURE, on a l'expression suivante :

- pour une signature donnée et une position donnée, on a une valeur et le coefficient correspondant

Le champ POSITION a été créé pour prévoir que les valeurs de la signature puissent être équivalente donc une valeur est associée à une position dans une signature. En ce qui concerne le champ COEFFICIENT, il permet de comparer deux images entre elles par le biais

de la signature. Cette comparaison est une des méthodes de recherche d'images par le contenu exprimées dans les besoins. Elle sera expliquée plus en détails dans la partie « Réalisation » et plus précisément dans la sous partie « Création du mini site ».

VIII- Réalisation

A- Développement du module Interface

Le développement de ce module se fait via la bibliothèque de classes Java Swing. Les outils choisis pour cette création se basent sur les fonctionnalités que l'interface doit avoir afin de reprendre parfaitement la plateforme SOX.

Voilà un aperçu de l'applet au moment de son lancement :

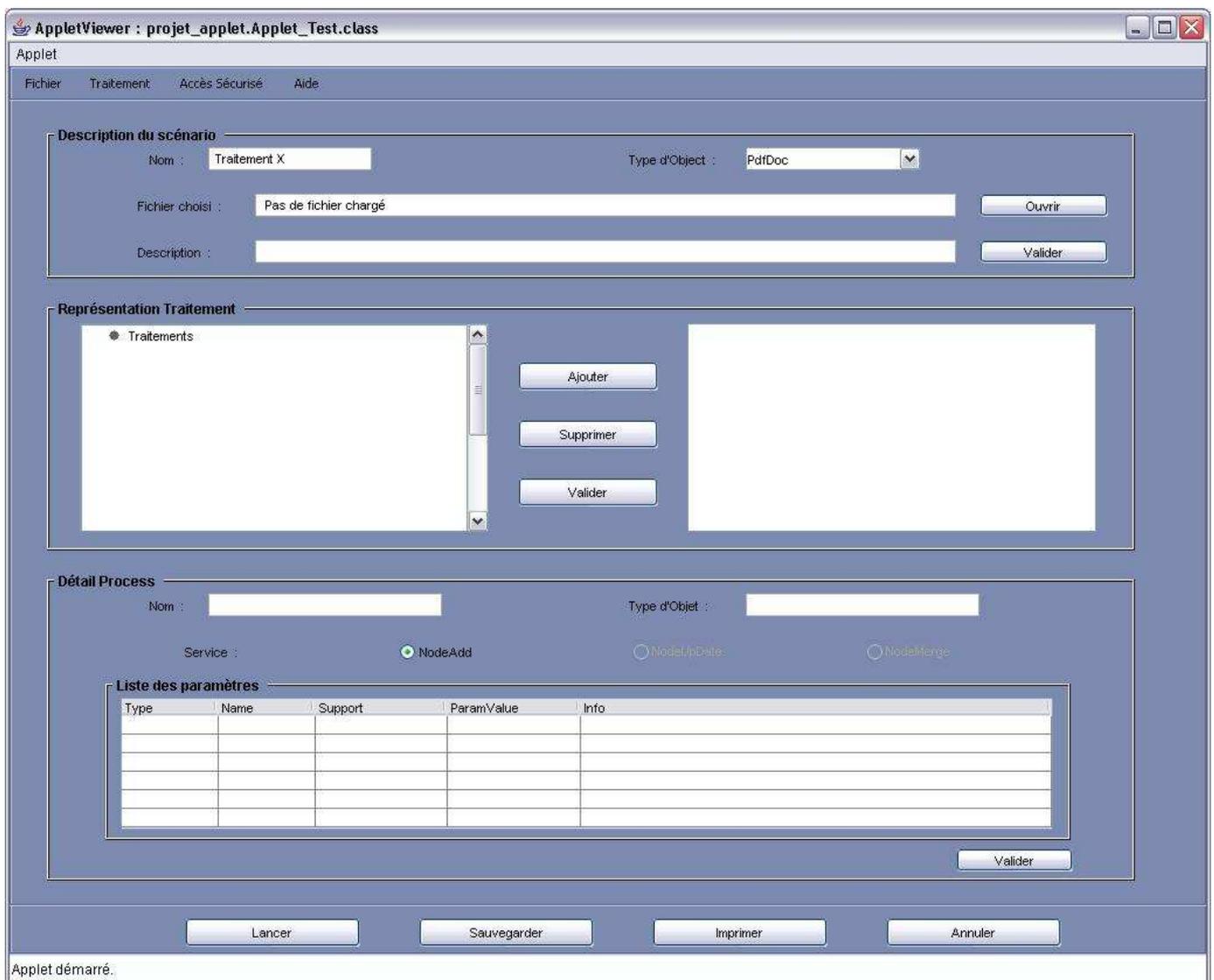


Figure 19 - Applet correspondant à l'interface SOX

Cette applet donne la possibilité à l'utilisateur de créer, modifier des scénarii permettant par la suite d'exécuter des traitements sur des documents. Il est également possible de sauvegarder ces scénarii mais aussi de les imprimer.

Lors de la création d'un scénario, l'utilisateur spécifie le nom du scénario, le type d'objet sur lequel va s'appliquer ce scénario, le document concerné et une description du scénario. Après avoir spécifié toutes ces données, l'utilisateur voit la JList, contenant tous les traitements, se modifié lui permettant de choisir le traitement à exécuter sur le document.

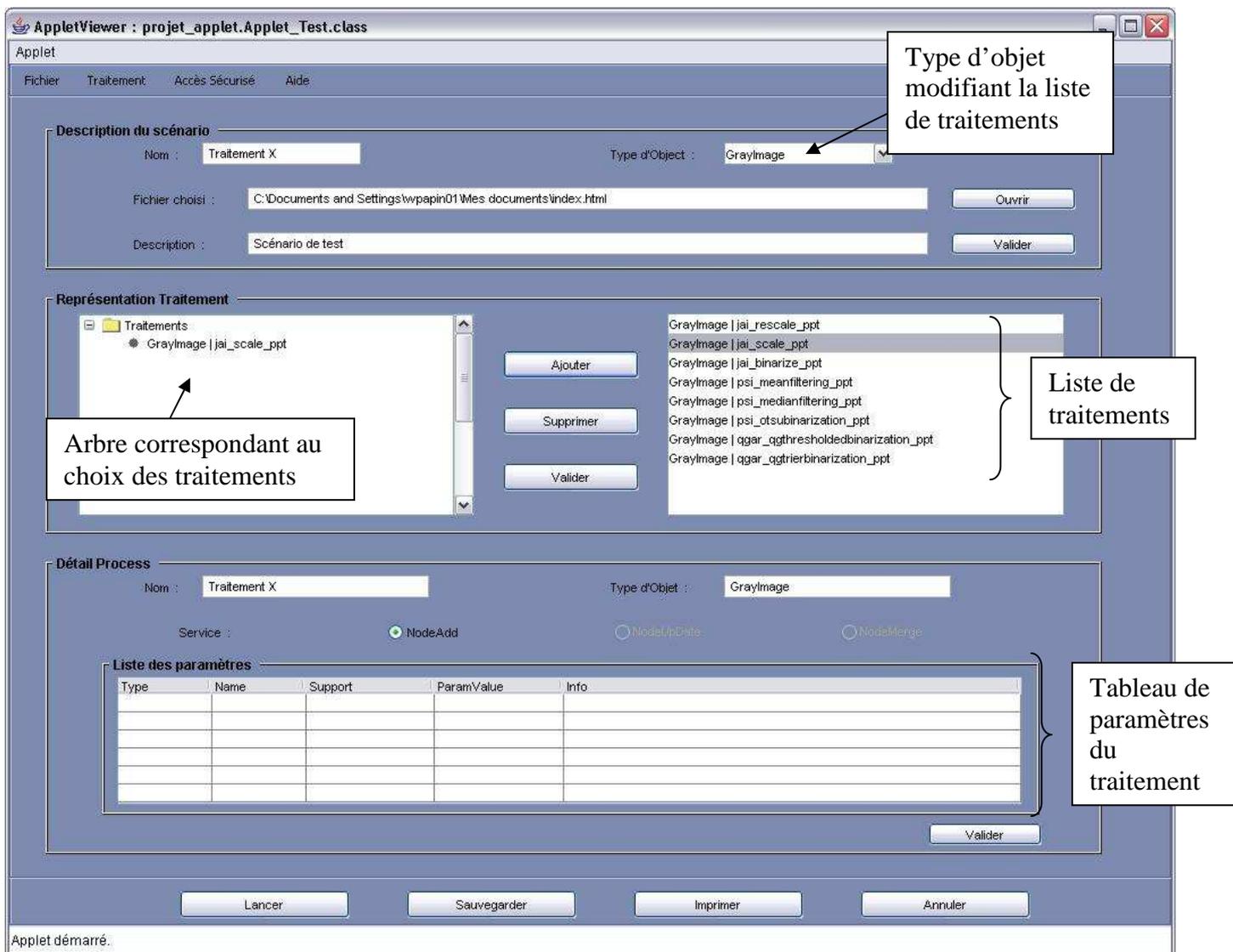


Figure 20 - Choix du traitement dans l'applet

Il est possible à l'utilisateur de l'applet de modifier les paramètres du traitement via la zone du bas de l'applet nommée « Détail Process ». Cette partie de l'applet est représentée par une JTable.

Concernant les autres fonctionnalités de l'interface, l'utilisateur peut modifier un scénario enregistré au préalable sur le serveur. Il lui suffit de cliquer sur le lien se trouvant dans la barre de menu.

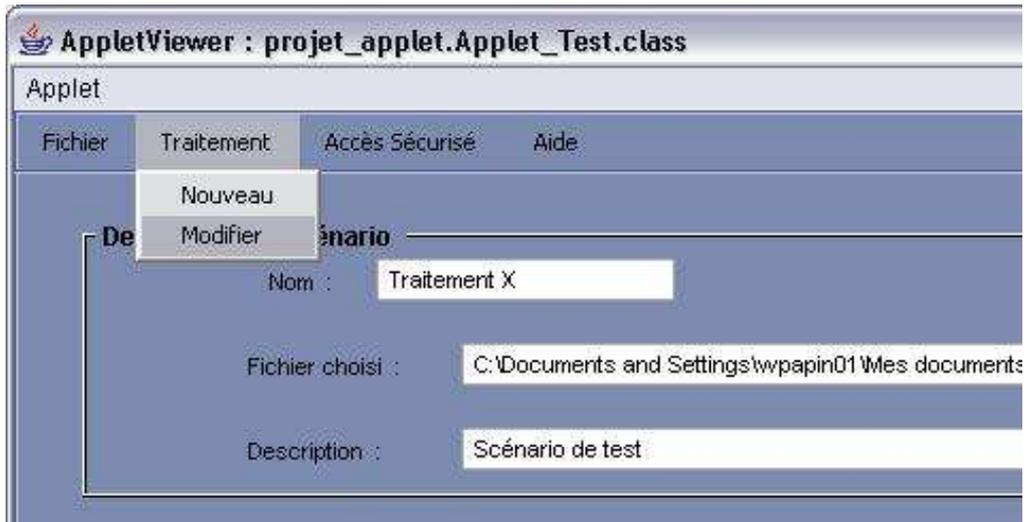


Figure 21 - Modification d'un scénario

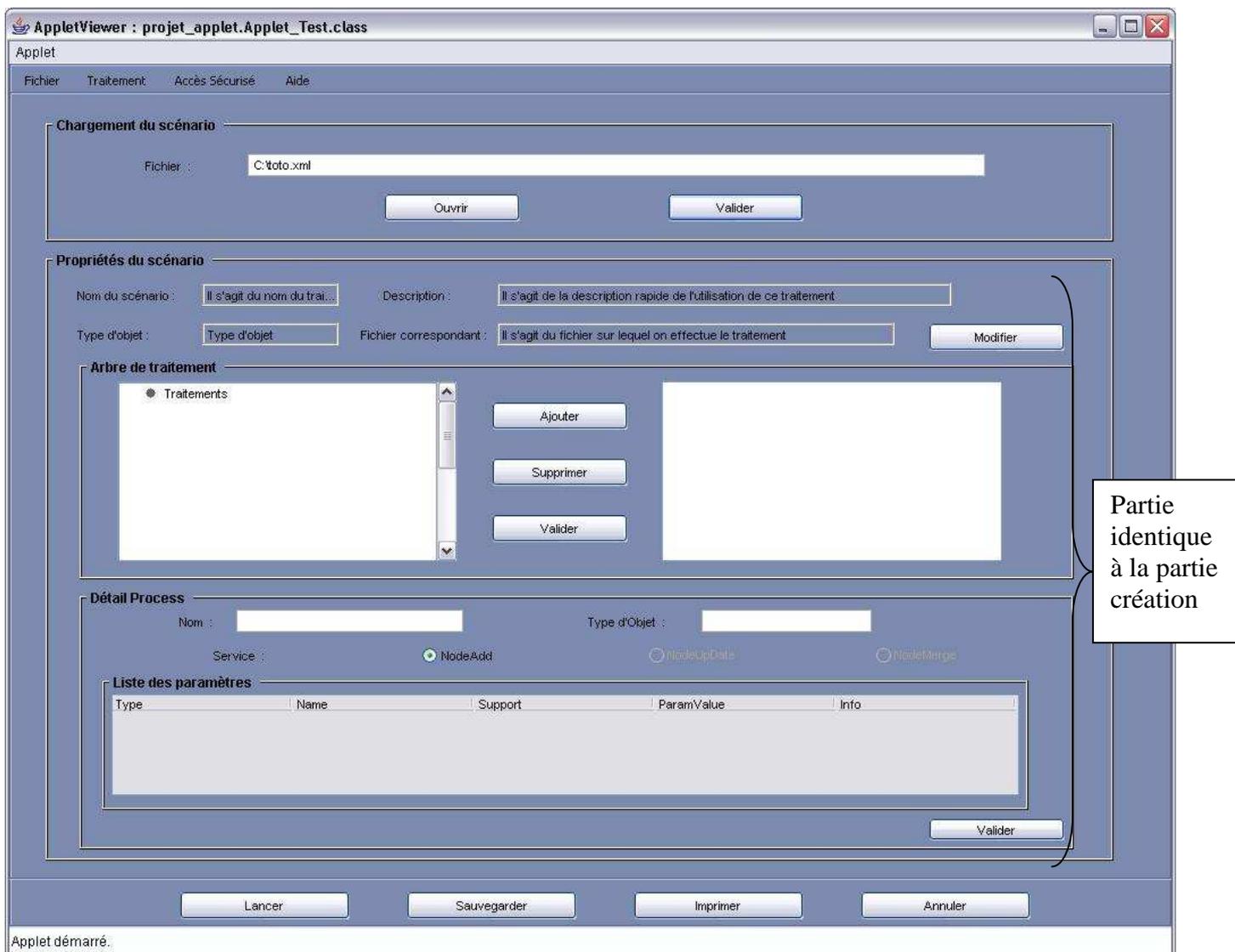


Figure 22 - Fenêtre de modification d'un scénario

L'applet donne la possibilité à l'utilisateur de sauvegarder, imprimer des scénarii. Concernant la sauvegarde, elle se fait par le biais d'un fichier XML contenant toutes les données présentes dans l'applet. Le fichier ainsi créé se nomme « save.xml » et se trouve à l'endroit choisi par l'utilisateur. Ensuite, pour l'impression d'un scénario, la technique utilisée passe par la technologie XML. En effet, l'utilisateur doit absolument enregistrer le scénario avant de pouvoir l'imprimer car l'applet imprime seulement un fichier XML. L'applet transforme ce fichier XML en graphique par le biais de la procédure « construitGraphe » transformant chaque ligne du fichier en graphique afin de faire l'impression.

Pour la partie Modification, l'utilisateur sélectionne le fichier de sauvegarde qu'il souhaite afin de le modifier. Et l'applet lit ce fichier XML et l'interprète via un « Parser » ou « Interpréteur ». La lecture se fait par le biais des balises XML du fichier que le « Parser » interprète afin de remettre les données dans les zones correspondantes de l'applet.

Voici un fichier de sauvegarde type :

```
<?xml version="1.0" encoding="ISO-8859-1" ?>
<savegarde>
  <description>Traitement X</description>
  <description_utilisation</description_utilisation>
  <fichier_correspondant>C:\Documents and Settings\wpapin01\Mes documents\index.html</fichier_correspondant>
  <type_fichier>PdfDoc</type_fichier>
  <traitements>

    <traitement1>
      <nom>PdfDoc | pdf_segment_ppt</nom>
      <parametres>

        <type>ParamIn</type>
        <name>Source</name>
        <support>StringValue</support>
        <paramvalue>@source</paramvalue>
        <info></info>

        <type>ParamIn</type>
        <name>ExtractImage</name>
        <support>String</support>
        <paramvalue>0</paramvalue>
        <info>if 1 extract the images</info>

        <type>ParamIn</type>
        <name>BigImages</name>
        <support>String</support>
        <paramvalue>1</paramvalue>
        <info>if 0 reduce the size of the image to fit original size</info>

        <type>ParamIn</type>
        <name>BigImages</name>
        <support>String</support>
        <paramvalue>1</paramvalue>
        <info>if 1 remove word textpieces (make the document lighter)</info>

      </parametres>
    </traitement1>
  </traitements>
</savegarde>
```

Figure 23 - Fichier XML de sauvegarde

La dernière partie de l'applet permet à un utilisateur possédant des droits d'administrateur de s'identifier dans la partie Accès Sécurisé afin d'ajouter un traitement à la

plateforme. Cette fonctionnalité permettant pour le moment de créer un fichier XML correspondant au traitement et de le voir s'ajouter à la plateforme logicielle.

Voici la partie d'identification :

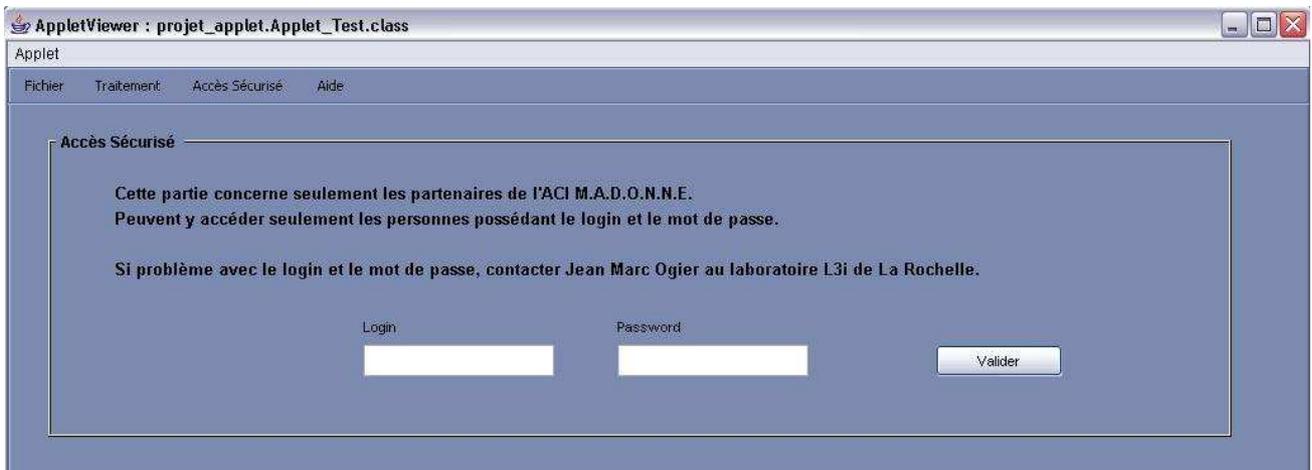


Figure 24 - Partie d'identification

Cette partie permet d'accéder à la zone d'ajout d'un traitement dont l'utilisateur doit spécifier le nom du traitement, le type d'objet concerné et bien évidemment l'adresse du serveur :

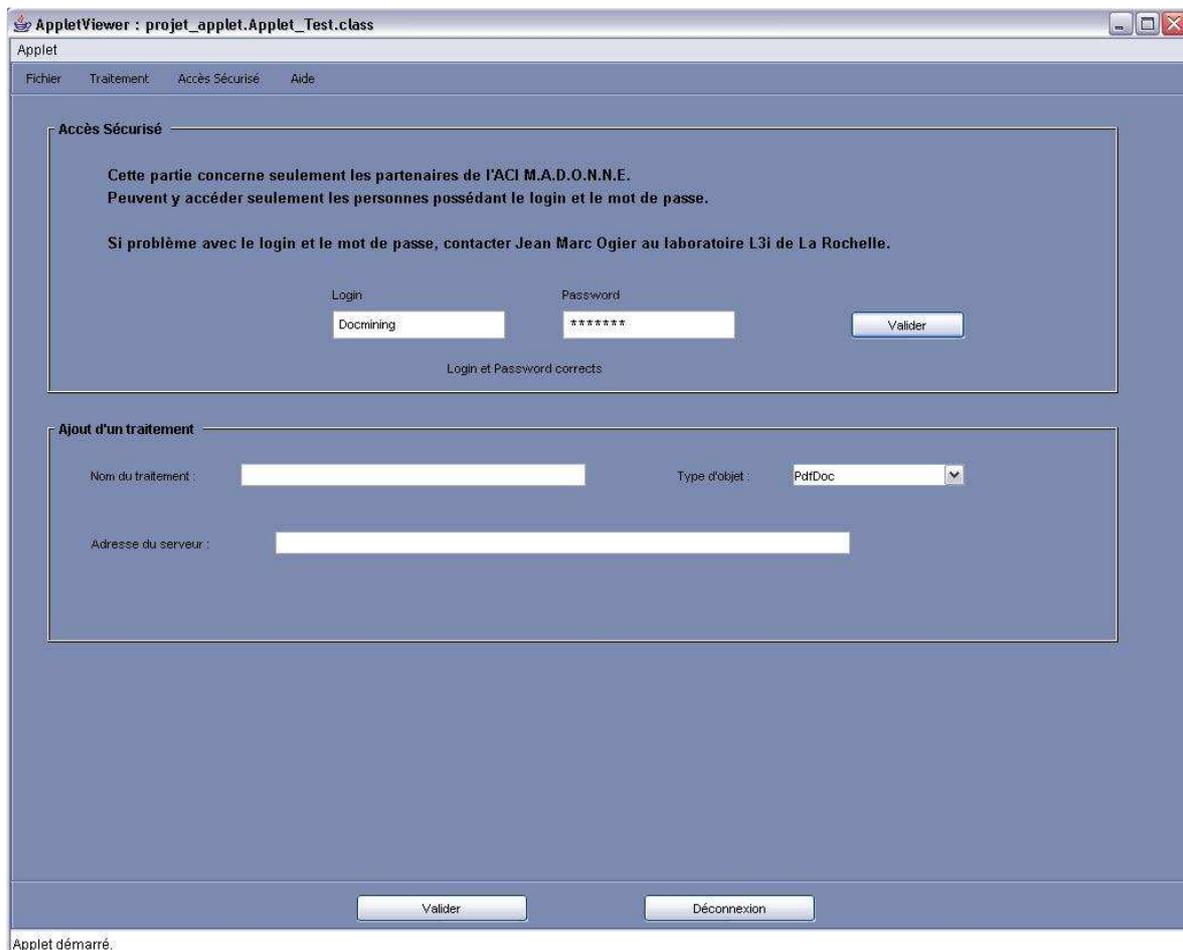


Figure 25 - Zone d'ajout

B- Intégration du module de traitement à l'interface

Après la création de la partie Interface de la plateforme, la partie Traitement de l'interface est à intégrer à cette applet.

C- Portage de la plateforme logicielle sur Internet

Comme le rapport l'explique dans la partie Conception, le portage de la plateforme se base sur une applet et sur une communication entre cette applet et le serveur par le biais d'une servlet.

L'utilisation de l'applet via un navigateur Internet se fait via l'archive jar de l'applet et par un fichier JSP dont voici le code :

Code:

```
<html>
<head>
  <title>
    Applet
  </title>
</head>
<body>
  <h2 align="center">Communication Applet - Servlet</h2><hr>
  <p align="center">
    <applet
      codebase = "."
      code    = "projet_applet.Applet_Test.class"
      archive = "Projet_Applet.jar"
      name   = "Applet_Test"
      width  = "1024"
      height = "768"
      hspace = "0"
      vspace = "0"
      align  = "middle">
    </applet>
  </p>
</body>
</html>
```

Lors du chargement de l'applet dans le navigateur Web, une fenêtre d'autorisation de chargement apparaît.

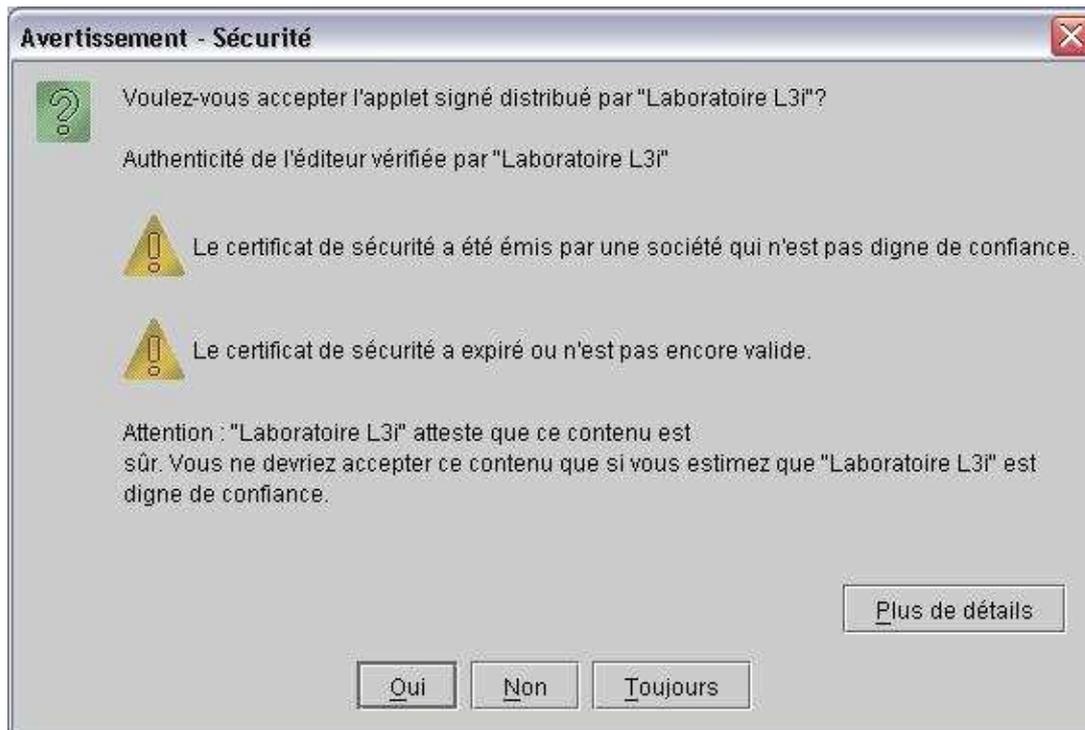


Figure 26 - Fenêtre de sécurité pour autorisation de l'applet

D- Création du moyen de stockage

La création de la base de données se fait par le biais de scripts SQL.

Ces scripts permettent de faciliter l'accès à la base de données mais aussi d'ajouter des champs qui seront utiles lors de l'utilisation du mini site de requêtes. Voici un des scripts de création de la table T_TYPE_IMAGE_TIM :

```
create table T_TYPE_IMAGE_TIM  
(  
    TIM_ID           INTEGER           NOT NULL,  
    TIM_FORMAT      CHAR (4)           NOT NULL,  
    TIM_LIBELLE     VARCHAR (32)       NOT NULL,  
    primary key (TIM_ID)  
)
```

Les autres scripts sont visibles dans les annexes.

E- Création du mini site

En plus de la création de la base de données d'images, cette seconde partie du sujet de stage comprenait la mise en place d'un mini site de requêtes permettant à un utilisateur d'effectuer des recherches d'images dans la base via des critères tels que la signature ou les méta-données, propre à chaque image.

Pour cela, une rapide étude fût faite afin de recenser les besoins des utilisateurs dans la recherche d'image. Il en est ressorti la liste suivante :

- recherche d'images via une signature de test;
- recherche d'images par méta-données via une liste ;
- affichage de toutes les images se trouvant sur le serveur pour une vue d'ensemble ;
- insertion rapide et facile d'images dans la base pour les administrateurs ;

Des screenshots du mini site seront présents dans la documentation technique présente en Annexes.

Cette partie va maintenant vous présenter les trois fonctionnalités de ce mini site.

1- Recherche d'image par critères

La recherche d'images se fait via plusieurs critères :

- Différence avec signature de référence
- Méta-données

a) Par signature

Il s'agit d'une suite de valeurs numériques correspondant à une image. Par exemple, une image A, suite à un traitement permettant de trouver sa signature numérique, va avoir comme signature les valeurs suivantes : 3,2 – 1,45 – 4,6. (Des propriétés du traitement différentes permettent d'avoir des valeurs différentes mais aussi un nombre de valeurs autre que trois).

Cette notion de signature est un peu équivalente à celle utilisée dans l'aviation pour les radars. Ces derniers, afin de détecter des avions ennemis et reconnaître les avions alliés utilisent la signature de l'appareil. Les avions alliés ont une signature connue du radar par le biais de sa base de données alors que l'avion ennemi n'en a pas ou pas de celles se trouvant dans la base.

Pour revenir au stage, la recherche d'images par ce critère oblige l'utilisateur à écrire dans la page correspondante une signature de référence permettant ensuite de trouver les images ayant une signature se rapprochant le plus de celle de référence :



Figure 27 - Recherche d'images par signature

Voilà le résultat de la recherche via une signature de référence de 3 valeurs :

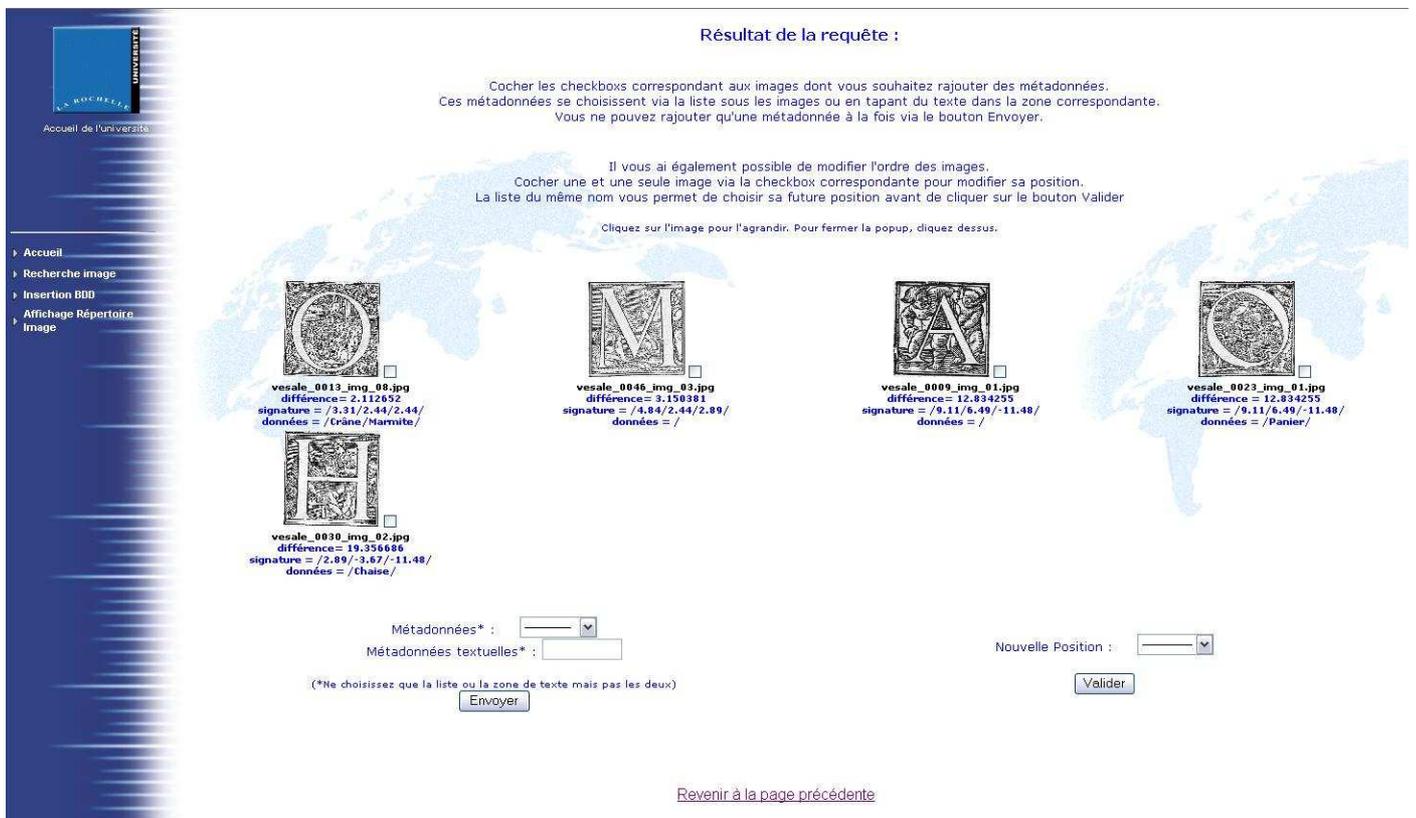


Figure 28 - Résultat d'une recherche par signature

Chaque image résultat est décrite par son nom, les méta-données lui correspondants, sa signature et la différence avec celle de référence. En réalité, l'affichage du résultat se fait dans un ordre précis (croissant) par rapport à la différence entre la signature de test et la signature de chaque image de la base (signature ayant le même nombre de valeur que celle de référence). Et la différence se calcule de manière simple par l'équation suivante :

$$différence = \sqrt{\left(\sum (x_i - y_i)^2\right)}$$

Les valeurs de i allant de 1 au nombre de valeurs de la signature. Au niveau de la base de données, la requête permettant d'effectuer ce calcul est très difficile car elle intègre des notions des modules mathématiques comme la somme au carré, la racine carré mais aussi la notion de boucle allant de i au nombre de valeurs de la signature. Cette notion m'étant totalement inconnue avant cette partie.

Voici la requête, séparée en plusieurs 'morceaux', permettant ce calcul :

```
$query = "SELECT i.IMG_ID, si.SIM_PATH, i.IMG_NOM_FICHER,
imgs.SIGNATURE_ID, SQRT (SUM(vs.COEFFICIENT*(POWER (v.VALEUR - ";

for ($i=1;$i<$nbchamp;$i++)
{
    $query = $query."IF(vs.POSITION - ".$i.", ";
}

$query = $query.$_POST['c'].$nbchamp];

$i = $nbchamp - 1;
while ($i != 0)
{
    $query = $query.",$_POST['c'].$i.");
    $i--;
}

$query = $query.",2))))";

$query = $query."FROM T_IMAGE_IMG i, T_VALEUR v,
T_IMAGE_SIGNATURE imgs, T_VALEUR_SIGNATURE vs,
T_STOCKAGE_IMAGE_SIM si
WHERE i.IMG_ID = imgs.IMG_ID
AND i.SIM_ID = si.SIM_ID
AND imgs.SIGNATURE_ID = vs.SIGNATURE_ID
AND vs.VALEUR_ID = v.VALEUR_ID
GROUP BY imgs.IMG_ID, imgs.SIGNATURE_ID
HAVING count(*) = $nbchamp
ORDER BY 5;";
```

La figure suivante illustre l'ensemble de ces descriptions ainsi que le résultat de la requête précédente.

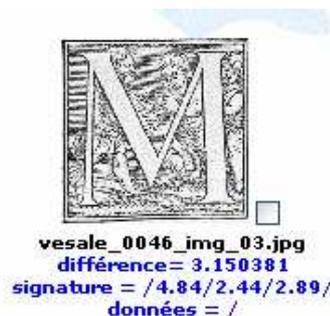


Figure 29 - Résultat d'une recherche par signature

Cet exemple met en évidence la présence d'une case à cocher à droite de l'image. Cette case donne une fonctionnalité en plus à ce mini site. En effet, l'utilisateur a la possibilité de modifier l'ordre des images s'il pense qu'une image ne devrait pas se trouver après ou avant une autre. Il suffit qu'il sélectionne l'image se trouvant mal placée et choisisse dans la liste la future position pour que l'ordre en soit modifié. Mais bien entendu, afin que cet ordre soit respecté la prochaine fois, des coefficients de pondération ont été rajouté à chaque valeur de chaque signature. L'utilisateur modifie donc les coefficients selon son choix.

Il est également possible, via ces cases à cocher, d'ajouter des méta-données à des images. L'utilisateur coche les images puis choisit une méta-donnée sachant qu'il a la possibilité de la faire par le biais d'une liste prédéfinie ou en écrivant directement le nom de la méta-donnée dans une zone de texte. Par la suite, cette méta-donnée est enregistrée dans la base de données via un lien avec l'image ou les images sélectionnées.

b) Par méta-données

La recherche par méta-données est beaucoup plus simple. Une image peut être composé de différents objets tels une chaise ou un livre. Ces objets que l'on appelle méta-données sont propres à chaque image mais un livre peut se trouver dans plusieurs images.

La requête exécutée pour la recherche par méta-données est la suivante :

```
$query = SELECT i.IMG_ID, si.SIM_PATH, i.IMG_NOM_FICHER  
FROM T_IMAGE_IMG i, T_STOCKAGE_IMAGE_SIM si, T_OBJET t,  
T_IMAGE_OBJET ti  
WHERE i.SIM_ID = si.SIM_ID  
AND i.IMG_ID = ti.IMG_ID  
AND ti.OBJET_ID = t.OBJET_ID  
AND t.OBJET_LIBELLE = '$objet' ;
```

Cette requête sélectionne les images via l'identifiant, le chemin d'accès à l'image sur le serveur et bien évidemment le nom de l'image selon celle ayant la variable **\$objet** comme méta-données.

Voici la présentation d'une image :



Figure 30 - Résultat d'une recherche par méta-données

Une fonctionnalité supplémentaire est rajoutée dans ce résultat, il s'agit de la possibilité de visualiser l'image dans sa taille réelle afin de mieux voir les méta-données directement sur l'image. Pour se faire, il suffit de cliquer sur l'image et une « popup » apparaît comme sur la figure suivante :

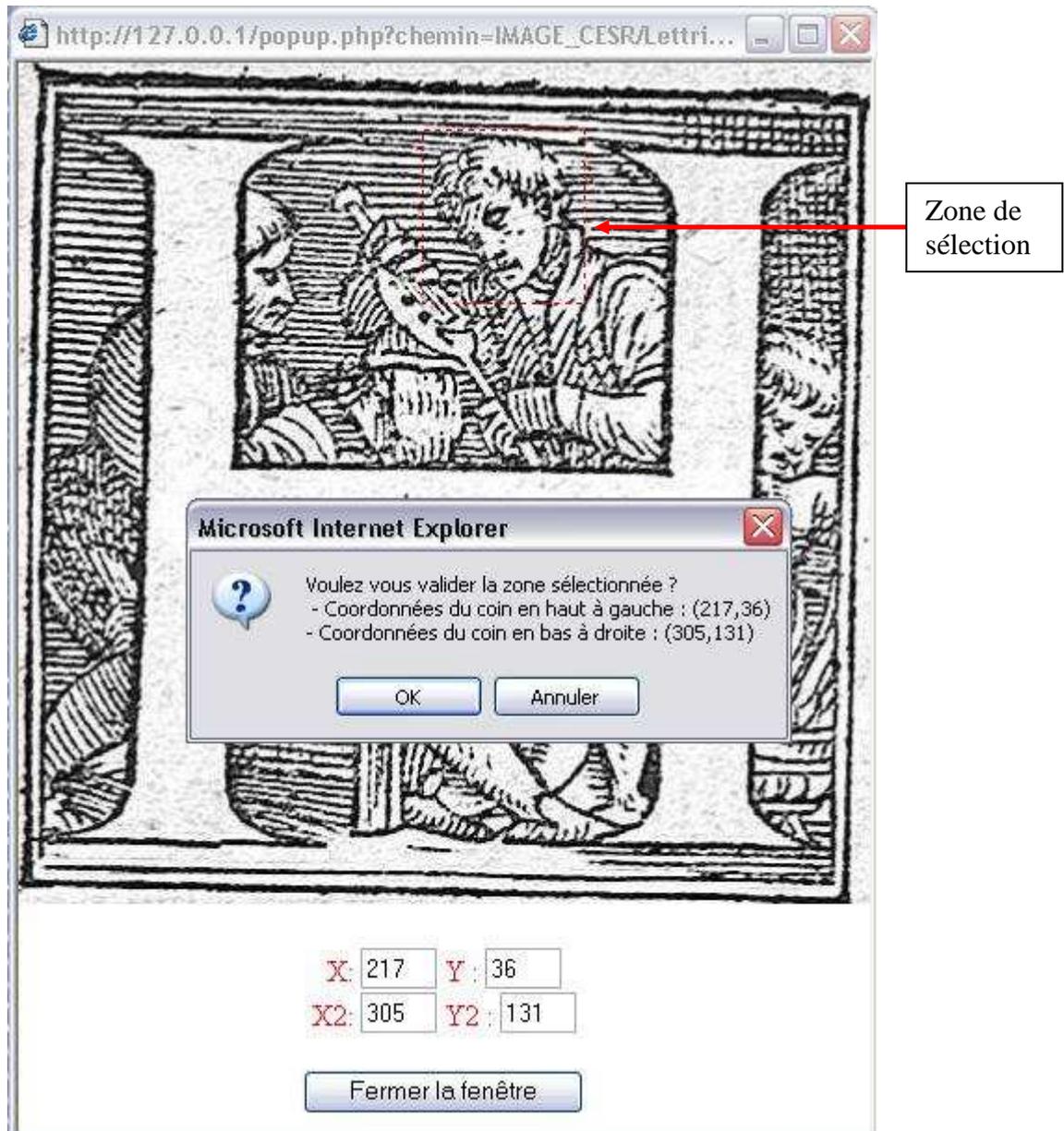


Figure 31 - Popup d'affichage de l'image

Cette « popup » donne également la possibilité à l'utilisateur de sélectionner une zone de l'image (visible par la zone pointillée rouge). Elle donnera la possibilité à l'utilisateur d'enregistrer cette petite image dans la base s'il le souhaite par le biais de la validation de cette petite fenêtre (visible sur l'image ci-dessus) récapitulant les coordonnées de la zone sélectionnée.

2- Affichage Répertoire

Cette partie n'est pas techniquement difficile mais elle a quand même une grande utilité. Elle va permettre de visualiser toutes les images d'un répertoire se trouvant sur le serveur comme on peut le faire avec un logiciel de photos par exemple.

(Voir la documentation technique dans les Annexes pour visualiser le résultat)

3- Insertion dans la base

Cette partie est réservée aux administrateurs du site. En effet, la page concernant l'insertion d'une image dans la base est protégé par un mot de passe permettant de sécuriser cette zone. Le mot de passe pouvant être piraté, il a également été crypté via un cryptage MD5 rendant tout accès à la page impossible sans le login et le mot de passe.

Ensuite, pour les possesseurs du mot de passe, ils peuvent, par le biais d'un formulaire, ajouter une image avec les données la caractérisant dans la base.

L'image suivante montre ce formulaire :

Insertion d'images dans la base

Nombre de champs de votre signature :

Nombre de métadonnées à rajouter :

signature :

métadonnées :

Nom fichier :

Description* :

(*par exemple : Lettrine R)

Taille :

Copyright :

Type d'image :

Figure 32 - Formulaire d'insertion d'images

La partie suivante de ce rapport fait une analyse des résultats obtenus et montre leurs adéquations par rapport aux objectifs initialement prévus.

IX- Objectifs fixés et résultats obtenus

Les objectifs principaux de ce stage étaient, premièrement, de porter la plateforme logicielle SOX sur Internet afin qu'elle soit utilisable à distance via un navigateur web. Deuxièmement, de créer une base de données d'images et un outil de recherche par critère de ces images.

Dans un premier temps, mes travaux se sont essentiellement basés sur le plateforme. Le problème survenu à mon arrivé avec l'absence du module Interface de la plateforme n'a pas aidé à la réussite totale des objectifs initiaux. La création de ce module a fait perdre beaucoup de temps. De plus, l'absence de documentation pour le projet Docmining (projet précédent et initiateur du projet sur lequel j'ai travaillé, le projet M.A.D.O.N.N.E), n'a pas aidé à la compréhension de celui-ci. Pour conclure sur cette partie, l'objectif initial était de porter la plateforme sur internet mais les résultats n'ont pas été aussi poussé car le module Interface a bien été refait avec de nouvelles fonctionnalités, l'intégration au module de Traitement a été réussi mais le portage via l'outil Servlet n'a pas encore abouti. De nombreuses erreurs sont encore présentes lors de l'utilisation de la plateforme. Le déploiement n'a donc pas pu être fait sur le serveur de test même si le travail a quand même été présenté à mon maître de stage. Je me suis accordé, après discussion avec lui, un certain temps de travail personnel pour essayer de régler certains problèmes essentiellement liés à l'utilisation des servlets. D'ici la présentation, certaines erreurs auront peut être été réglé.

Concernant la deuxième partie du stage, le travail a pu être réalisé dans les délais. L'objectif était de créer une base de données d'images via l'outil adéquate et de réaliser un outil de recherche d'images par le biais de critères. La perte de temps découlant de la première partie du stage n'a pas aidé pour cette partie donc comme il n'y avait pas d'obligations quant aux moyens à utiliser pour la réalisation de ces outils, les outils les plus simples d'utilisation ont été choisis (MySQL pour la base de données et un outil internet pour la recherche multicritères).

Au final, la création de la base s'est déroulée sans aucun problème. Comme l'outil choisi ne permettait pas de gérer les images elles-mêmes et que cette gestion n'est pas toujours la meilleure, j'ai choisi de gérer le chemin d'accès aux images au niveau de la base. Ensuite, pour l'outil de recherche, le travail via des pages internet fût choisi et toutes les fonctionnalités demandées ont été ajouté.

Le déploiement a donc pu être réalisé sur le serveur internet de test moyennant quelques petits problèmes de compatibilité. En effet, le code PHP créé ne correspondait pas du tout au moteur PHP du serveur. Le problème était dû aux balises, différents selon les moteurs. Cette partie là fût validée par mon maître de stage et l'outil a finalement pu être intégré au site du projet MADONNE et différents tests de fonctionnement ont été réalisé afin de régler les dernières petites erreurs de présentation.

Pour conclure sur les résultats obtenus, ils ont été en deçà des objectifs initiaux. Essentiellement dû à la perte de temps passé sur la création du module Interface et mes problèmes de compréhension du projet Docmining et de l'utilisation des servlets.

Ces résultats permettent quand même de penser aux perspectives. L'outil de recherche d'images pourra par la suite être intégré à un module de traitement calculant la signature de certaines images et insérant automatiquement après calcul, la signature dans la base.

Pour la plateforme, les perspectives sont bien évidemment de régler les problèmes liés au portage (erreurs de servlets). Par la suite, de permettre à tous les administrateurs d'envoyer les codes des nouveaux traitements à ajouter afin de les utiliser automatiquement sur la plateforme après l'ajout.

Les parties suivantes de ce rapport décriront les acquis apportés par le stage en terme de compétences techniques et les aspects humains du stage.

X- Acquis techniques

Le stage qui m'a été proposé se sépare en deux parties bien distinctes. La première partie avait pour but de porter une plateforme logicielle en Servlet afin d'être utilisé à distance via un navigateur Web. Le déroulement du stage a amené à refaire la plateforme logicielle via le langage Java et ses bibliothèques d'interface. Ce travail ne m'a pas permis d'acquérir de nouvelles compétences, s'agissant d'un langage que je maîtrisais déjà un peu. Cependant, certaines fonctionnalités de la plateforme m'ont fait utiliser des outils du langage Java inconnus pour moi précédemment comme un Parser XML (outil permettant d'interpréter n'importe quel fichier XML par le biais de ces balises).

Par contre, le travail réalisé pour le portage de la plateforme sur Internet fût un acquis conséquent quant aux compétences techniques. L'utilisation des servlets afin de faire communiquer l'applet et le serveur via des requêtes HTTP a été un outil très intéressant à manipuler et a ouvert bien évidemment des portes vers des outils plus complexes comme Struts. (Outil non intégré dans le projet par faute de temps mais dont l'analyse a été faite afin de voir les différences avec les Servlets).

Concernant la seconde partie du stage sur la création de la base de données, les compétences techniques requises n'étaient pas très importantes. Cette création s'est faite via l'outil MySQL, outil déjà utilisé au préalable et assez simple à manipuler.

Pour finir sur cette partie, la création du mini site m'a permis d'approfondir mes connaissances dans le langage Javascript, langage que j'avais utilisé peu souvent. Les fonctionnalités de recherche d'images via certains critères se faisaient par le biais de requêtes. Certaines d'entre elles ont été plus difficiles à créer car faisant intervenir des notions de mathématiques et des boucles de lecture. Ce sont les seuls acquis techniques concernant cette partie.

J'ai pu sinon acquérir quelques techniques de traitements d'images permettant de transformer par exemple un fichier PDF en une image ou de décomposer une image contenant du texte en un arbre avec toutes les phrases du texte voir même toutes les lettres. Ce sont des techniques qui ne me serviront peut être jamais mais qui viennent enrichir mes compétences en général.

La partie suivante fait une analyse des aspects humains du stage.

XI- Bilan humain

Ce stage se positionnant dans un projet de grande envergure qu'est le projet M.A.D.O.N.N.E, intégrer une équipe de chercheurs afin de réaliser des outils logicielles m'a déjà apporté un bilan humain positif. Ayant été en plus étudiant dans l'Université où se situe le laboratoire, mon intégration a été facilitée par le biais de mes connaissances du monde universitaire et des chercheurs y travaillant.

Concernant l'intégration dans l'équipe M.A.D.O.N.N.E, elle fût facilité par la fonction de mon maître de stage dans ce consortium : coordinateur. De plus, je fus présenté à l'ensemble des membres du consortium lors d'un état de l'art des différents sous-projets à Tours. J'ai pu d'ailleurs y effectuer une présentation de mes travaux devant l'ensemble des partenaires. Ces échanges avec des chercheurs venant de différentes régions de la France m'ont permis d'acquérir des connaissances dans les domaines de l'Imagerie mais m'ont également permis de régler quelques problèmes liés à mon stage via des discussions avec certaines chercheurs.

Par contre, cet état de l'art s'étant déroulé au cours du mois de Mai, j'ai dû planifier tout mon travail en fonction de cette présentation et réaliser les parties les plus importantes visuellement afin de montrer à la communauté M.A.D.O.N.N.E l'avancement des travaux du point de vue Génie Logiciel.

Au cours de cette présentation, de nombreux problèmes informatiques sont intervenus, la ralentissant très sensiblement et ne me permettant pas de montrer tous les travaux que je souhaitais exposer. Ces complications m'ayant permis malgré tout de progresser du point de vue expression orale. J'en ai tiré de nombreux enseignements. Le principal étant de prévoir tout en amont pour ne pas rencontrer de problèmes en aval.

En plus de tout le travail effectué durant ce stage, j'ai eu l'occasion, en parallèle, d'apporter certaines de mes connaissances à des étudiants de l'IUP Génie Informatique de La Rochelle par le biais de cours de travaux pratiques que j'ai dispensé à ces étudiants dans le domaine de la Base de Données. Au cours de ces TPs, j'ai pu me rendre compte du gros travail à réaliser si l'on veut un bon déroulement de séances mais j'ai également profité de ces séances pour fortifier mes bases dans ce domaine.

La dernière partie de ce rapport conclura sur ce stage, du déroulement technique, humain pour finir par les perspectives futures.

Conclusion

Les membres du consortium M.A.D.O.N.N.E (projet monté par différents chercheurs du laboratoire L3i et de laboratoires partenaires répondant à un appel à projet du Ministère de la Recherche et du CNRS) ont développé une plateforme logicielle, nommée SOX, permettant d'intégrer des outils hétérogènes de traitement. L'objectif de mon stage, intégré dans ce projet, est de porter cette plateforme sur internet afin d'être utilisée à distance. Dans un second temps, il s'agit de créer une base de données d'images dans un but d'organisation et de préservation des données patrimoniales.

Pour la première partie du stage, ma solution s'appuie sur un outil permettant de faire communiquer le module Interface de la plateforme et le module Traitement se trouvant sur le serveur internet, les servlets. Pour la seconde partie, la création de la base de données d'images réside dans le fait de gérer le chemin vers ces images et non les images elles-mêmes. Toutes les étapes d'un projet sont présentes dans ce stage : de l'étude de l'existant à l'analyse fonctionnelle, en passant par la conception et la réalisation et pour finir sur les tests.

Les résultats obtenus ne sont pas en total conformité avec les attentes du consortium M.A.D.O.N.N.E. La partie de création de la base de données ainsi que l'outil permettant de faire des recherches par critères dans cette base atteint tous les objectifs initiaux. Par contre, la partie concernant le portage de la plateforme sur internet s'avère ralenti par le fait que le module Interface de cette plateforme a dû être totalement refait à cause d'un problème interne au laboratoire mais aussi par le fait de l'absence totale de documentation technique sur la partie Traitement (partie contenue dans un précédent projet, le projet Docmining). Actuellement, le portage génère encore de nombreuses erreurs liées à la communication entre le serveur et l'applet. Le module de création et d'utilisation de la base est par contre d'ores et déjà validé par Mr Ogier.

J'ai eu beaucoup de plaisir à travailler au laboratoire L3i durant ce stage. L'ambiance de travail et l'esprit de confiance qu'il y avait entre Mr Ogier et moi-même m'a permis de travailler dans de bonnes conditions même si le stage n'a pu être totalement mené à termes au niveau des objectifs initiaux.

Dans la continuité de la formation dans laquelle je suis intégré, une partie Gestion de projet a été réalisée. Dans ce cadre, une estimation des risques a été réalisée. En règle générale, certains risques sont prévus et suivis mais d'autres peuvent apparaître tout au long du projet. Dans mon cas, il s'avère que certains risques non prévus sont apparus, ralentissant notablement le stage. J'ai essayé de les gérer de la meilleure manière possible en suivant les conseils de nos professeurs durant ma formation.

Ce projet constitue pour le consortium M.A.D.O.N.N.E une petite étape dans la valorisation du patrimoine et dans la promotion de ces travaux vers l'extérieur. Ce développement peut être applicable à d'autres travaux des partenaires du consortium. A terme, ces derniers souhaitent développer ce type d'outil dans le but final de valoriser les activités de numérisation du patrimoine français et européen.

Table des légendes

Figure 1 - Le personnel du L3i en 2003	8
Figure 2 - Le budget du L3i en 2002.....	8
Figure 3 - Plateforme SOX.....	15
Figure 4 - Résultat de l'interprétation du logiciel SOX.....	16
Figure 5 - Menu de traitements possibles sur un objet.....	16
Figure 6 - Diagramme de planification	18
Figure 7 - Diagramme de GANTT	19
Figure 8 - Feuille de risque	20
Figure 9 - Feuille de suivi des risques.....	21
Figure 10 - Infrastructure matérielle et logicielle.....	22
Figure 11 - Nouvelle infrastructure avec prise en compte du travail réalisé.....	24
Figure 12 - Fenêtre de visualisation XMillum	27
Figure 13 - Sélection du traitement à exécuter.....	27
Figure 14 - Résultat du traitement.....	28
Figure 15 - Fenêtre de paramètre	28
Figure 16 - Fonctionnement d'une servlet.....	31
Figure 17 - Fonctionnement de Struts	32
Figure 18 - Schéma de la base de données.....	35
Figure 19 - Applet correspondant à l'interface SOX	37
Figure 20 - Choix du traitement dans l'applet	38
Figure 21 - Modification d'un scénario	39
Figure 22 - Fenêtre de modification d'un scénario.....	39
Figure 23 - Fichier XML de sauvegarde	40
Figure 24 - Partie d'identification.....	41
Figure 25 - Zone d'ajout	41
Figure 26 - Fenêtre de sécurité pour autorisation de l'applet	43
Figure 27 - Recherche d'images par signature	46
Figure 28 - Résultat d'une recherche par signature	46
Figure 29 - Résultat d'une recherche par signature	47
Figure 30 - Résultat d'une recherche par méta-données.....	48
Figure 31 - Popup d'affichage de l'image.....	49
Figure 32 - Formulaire d'insertion d'images.....	50