

LOG 792	DOCUMENT NO. 1.0	DATE 2012-09-24	VERSION 1.0
TITRE		PAGE	PAGES
EndoMine - Rapport d'étape		1	38

# Département de génie logiciel et des TI

# RAPPORT D'ÉTAPE

LOG 792. Projet de fin d'études Département de génie logiciel et des TI

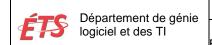
#### **EndoMine**

Projet de développement d'outils de forage de données de résultats de tests patients - endocrinologie, métabolisme et épidémologie clinique du JGH

Auteurs
Anton Zakharov
ZAKA12038406
David Lauzon
LAUD01028300

Professeur superviseur Alain April

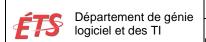
Date 29 Octobre 2012



COURS LOG 792	DOCUMENT NO. 1.0	DATE 2012-09-24	version 1.0
TITRE		PAGE	PAGES
EndoMine - Rapport d'étape		2	38

Suivi des changements \*A – Ajouté M – Modifié S – Supprimé

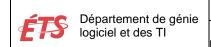
NUMÉRO DE VERSION	<b>DATE</b> aaaa/mm/jj	NUMÉRO DE FIGURE, TABLE OU SECTION	A* M S	BRÈVE DESCRIPTION DU CHANGEMENT	NUMÉRO DE DEMANDE CHANGEMENT
1.0	2012/10/28		А	Document Complete	



COURS LOG 792	DOCUMENT NO. 1.0	DATE 2012-09-24	version 1.0
TITRE		PAGE	PAGES
EndoMine - Rapport d'étape		3	38

# Table des matières

1. Problématique et contexte	
2. Objectifs du projet	
3. Méthodologie	
4. Sommaire des travaux réalisés et recommandations	
4.1 Sommaire des travaux réalisés	
4.1.1 Définition des besoins du client	
4.1.2 Exploration des différents technologies	
4.1.3 Test Effectués sur différentes technologies	
4.2 Recommandations	
5. Livrables et planification	
5.1 Description des artéfacts	
5.2 Planification	
6. Risques	13
7. Références consultées	15
8. Table des matières du rapport	15
Annexe A : Plan de travail révisé	17
Annexe B : RÉFÉRENCES	21
Annexe C : Installation de Hadoop et Hbase (dans linux)	30
Annexe D : Document de vision	38



COURS LOG 792	DOCUMENT NO. 1.0	DATE 2012-09-24	VERSION 1.0
TITRE		PAGE	PAGES
EndoMine - Rapport d'étape		4	38

#### 1. PROBLÉMATIQUE ET CONTEXTE

Actuellement le département de diagnostic médical du JGH possède une base de données de production Oracle contenant toutes les données diagnostics. Pour des raisons de sécurité des données et de performance, une base de données Access a été créée et elle est actuellement utilisée pour effectuer des extractions et des fonctions de « data mining ». Cette approche comporte des limites et le Dr. Eintrachtt aimerait avoir des propositions de solutions pour l'avenir d'un environnement de « data mining » à grande échelle. Parallèlement, le département d'endocrinologie du JGH aimerait avoir un outil de « data mining » pour effectuer des recherches spécifiques à son domaine d'expertise.

#### 2. OBJECTIFS DU PROJET

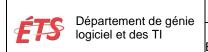
Notre projet contient plusieurs aspects :

- Adapter le scope du projet dépendant l'impact qu'il aurait sur différents intervenants.
- La mise en place d'une base de données exploitable par des outils de « data mining » qui remplacerait la base de données Access;
- Le choix ou le développement d'outils (à base de logiciels libres) de « data mining » afin de permettre l'étude des données d'une manière interactive. Nous avons décidé d'explorer les différentes technologies de l'écosystème «Hadoop » étant donnée leurs coût de mise à l'échelle peu élevée. Nous allons aussi explorer quelques solutions propriétaires.
- Reproduction de la B.D. de production Oracle actuelle dans l'environnement de « data mining ». Synchronisation entre la B.D. de production et l'environnement de « data mining » et de la base de données de staging.
- Création d'une interface permettant de : 1) faire la recherche dans l'environnement de « data mining » et 2) En affichage graphique des données (style tableau ou visuel).

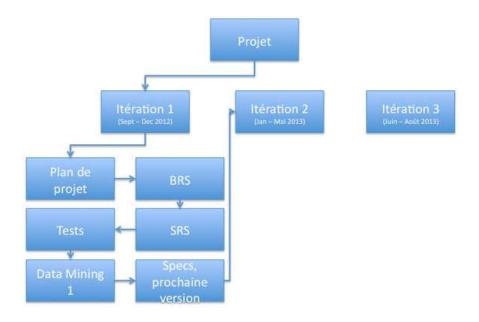
#### 3. MÉTHODOLOGIE

Vu la portée et la nature du projet l'approche itérative par phase (création, élaboration, construction, transition) du processus unifié (U.P.) sera utilisée. Les itérations et tâches exactes seront décrites dans l'Annexe A. Par contre certaines parties (comme par exemple reproduction de la BD de production) n'ont pas une définition exacte du temps, car nous n'avons aucune idée pour le moment quelle technologie nous allons utiliser. Ces questions seront résolues au fur et à mesure. Chaque itération prendrait 2 semaines (environ 15h de travail par personne par semaine). Un bilan serait fait à la fin de chaque évaluation, et les tâches seront réévaluées selon le temps disponible et l'évaluation des besoins du client.

Voici le schéma des livrable proposés :



COURS	DOCUMENT NO.	DATE	VERSION
LOG 792	1.0	2012-09-24	1.0
TITRE		PAGE	PAGES
EndoMine - Rapport d'étape		5	38



# 4. SOMMAIRE DES TRAVAUX RÉALISÉS ET RECOMMANDATIONS

#### 4.1 Sommaire des travaux réalisés

Pour la facile de présentation du projet nous allons diviser cette section en 3 parties : Définition des besoins du client, exploration des différents technologies et test effectuées sur différentes technologies.

Après nous allons décrire les installations effectuées.

#### 4.1.1 Définition des besoins du client

Étant donné que le projet est réalisé pour un client réel, la création d'un document de vision décrivant les différents intervenants et leurs besoins est essentielle. Le document de vision complet avait été ajouté à l'annexe D.

Les besoins des clients pouvant être divisé en 3 projets avec une architecture différente, nous avons décidé de se concentrer sur le projet répondant au Dr. Eintrachtt. Ce dernier voudrais pouvoir faire de recherches rapidement sur la base de donnée avec un générateur de requêtes de type Ms Access. Voici un exemple :



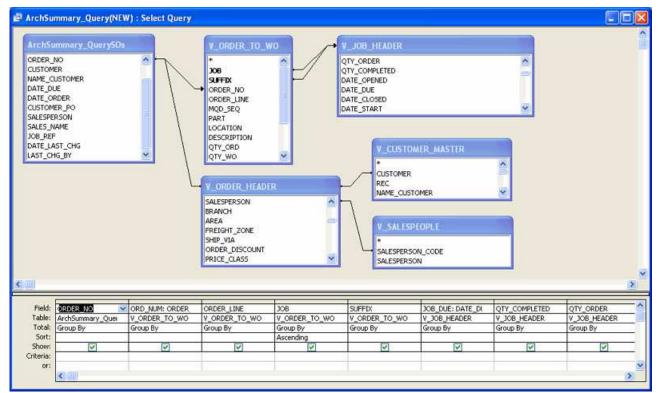
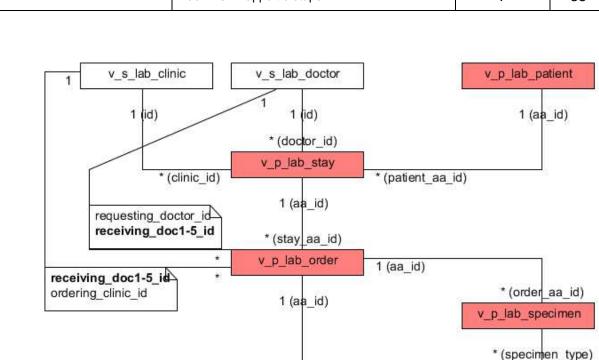


Figure 1. Générateur de requêtes MS Access.

Le nombre de tables dans la base de données est énorme. Comme nous ne faisons qu'un prototype, importer toutes les tables avec toutes les relations dépasse le scope de ce projet. Nous avons donc extrait et formaté les tables que le Dr. Einthracht allait utiliser.



\* (group\_test\_id)

\* (group\_test\_id)

\* (group\_test\_id)

\* (specimen\_type)

\* test\_id reflex\_test\_id

\* (order\_aa\_id)

\* (specimen\_type)

\* test\_id reflex\_test\_id

\* (specimen\_type)

\* (specim

Figure 2. Relations entre les tables principales.

Les tables commençant par v\_p contiennent les données particulières, alors que les v\_p servent à décrire l'information dans les tables v\_p. Par exemple, v\_s\_lab\_test contient le nom du test et les caractéristiques propres servant à classifier un test par rapport aux autres. V\_p\_lab\_test\_result contient des données comme les dates des tests réalisés, ainsi que les valeurs reliées aux tests.

Voici la description des tables :

- **V\_p\_lab\_patient** : est une table vide. Une clé anonymisée est utilisée pour rendre chaque patient unique.
- V p lab stay : contient les séjours de patients dans l'hôpital.
- V\_p\_lab\_order : contient la commande de tests effectués durant un séjour particulier.
- V\_p\_lab\_test\_result : contient les résultats des tests.

<b>∠</b> Département de génie	LOG 792	DOCUMENT NO. 1.0	DATE 2012-09-24	VERSION 1.0	
ER	logiciel et des TI	TITRE EndoMine - Rappor	t d'étape	PAGE 8	PAGES 38

- V\_p\_lab\_specimen : contient les spécimens physiques pris pour faire des tests.
- V\_s\_lab\_test\_group : décrit les composantes d'un test particulier.
- Les tables restantes n'ont pas besoins d'être décrits pour le moment.

Une fois les tables et les relations particulières effectuées, nous avons pris un soin de définir la taille de chaque table à partir d'un échantillon (de 1 jour de données). Voici les résultats obtenus :

		vp_lab			vp_lab_te	vs_la b_spe	vs_la	vs_la b_test _grou	
Data Type	Table	_stay	order	pecimen	st_result	cimen	b_test	p	TOTAL
		34	75	19	62	11	145	47	393
	D (								
Number(5)	Bytes 2	0	1	0	0	0	4	1	6
Number(10)									
	4	7	9	7	7	0	10	1	41
Number(14)	6	2	2	2	2	1	1	0	10
double	0	0	0	0	0	0	1.1	0	0
double	8	0	0	0	0	0	11	0	11 0
datetime	3	2	2	2	2	0	0	0	8
datetime	3	2				0	0	0	0
char(1)	1	1	7	0	1	0	7	1	17
varchar(1)	1	8	33	4	40	5	56	0	146
varchar(2)	2	0	1	0	1	0	0	0	2
varchar(3)	3	1	0	1	1	1	2	0	6
varchar(4)	4	0	0	0	0	0	1	0	1
varchar(5)	5	8	18	3	7	0	13	41	90
varchar(7)	7	1	0	0	0	0	0	0	1
varchar(10)	10	0	0	0	0	1	0	0	1
varchar(11)	11	3	2	0	0	0	28	0	33
varchar(15)	15	0	0	0	1	1	3	0	5
varchar(23)	23	1	0	0	0	1	0	1	3
varchar(30)	30	0	0	0	0	0	1	0	1
varchar(39)	39	0	0	0	0	1	0	0	1
varchar(59)	59	0	0	0	0	0	2	2	4
varchar(79)	79	0	0	0	0	0	4	0	4
varchar(239)	239	0	0	0	0	0	1	0	1
mediumtext	1024	0	0	0	0	0	1	0	1
									0
Fixed Size		46	56	46	46	6	142	6	348

		COURS	DOCUMENT NO.	DATE	VERSION
Département de génie logiciel et des TI	LOG 792	1.0	2012-09-24	1.0	
	TITRE		PAGE	PAGES	
	EndoMine - Rappor	t d'étape	9	38	

Variable Size	115	154	22	96	95	2218	347	3047
Max Size per								
row	161	210	68	142	101	2360	353	3395
(magnitude)	1	2	4	36				
per hour	84	140	367	3,048				3,639
per day	2,014	3,366	8,810	73,150	741	4,321	1,243	87,340
per week	14,098	23,562	61,670	512,050				611,380
per month	60,420	100,980	264,300	2,194,500				2,620,200
								31,879,10
per year	735,110	1,228,590	3,215,650	26,699,750				0

Table 1. Row count statistiques

Comme on peut voir la majorité des tables v\_p ont au moins 1 million de rangées par année. Faire les jointures sur des rangées de cette taille serait très difficile. Vue la taille des données, il va devenir très compliqué d'effectuer des jointures entre les grosses tables.

### 4.1.2 Exploration des différents technologies

Étant donné que le problème d'une grande quantité des données n'as pas de solution fixe, nous allons décrire les différentes méthodes analysées. Les bases de données Relationnelles, NoSQL Open Source et technologies propriétaires seront décrites.

#### 4.1.2.1 Base de données relationnelles

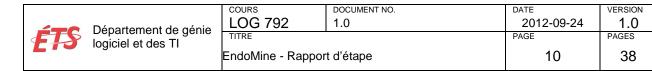
#### Oracle:

Avant de réinventer la roue, nous avons pensé à utiliser les bases de données existantes. En théorie la base de donnée Oracle, présentement en production devrait permettre de faire des requêtes nécessaires à Dr. Eintracht. Il nous faudrait les données au complet, ainsi que des requêtes spécifiques du Dr. Eintracht pour pouvoir le tester.

Cette solution va couter une licence Oracle et ne pourrait pas vraiment être mise à niveau si le nombre de données augmente.

# 4.1.2.2 Base de données NoSQL Open Source

Les bases de données NoSQL dont on va présenter vont être bâties en utilisant le système de gestion de fichier distribué «Hadoop». Présenter l'architecture de «Hadoop» vas au delà du scope du présente document, mais si vous êtes intéressé voici un lien qui la décrit assez sommairement : <a href="http://en.wikipedia.org/wiki/Apache\_Hadoop">http://en.wikipedia.org/wiki/Apache\_Hadoop</a> .



De base les données sont stockées dans des fichiers qui sont automatiquement répliqués sur plusieurs disques. MapReduce est un *framework* utilisé pour créer des tâches fonctionnent en parallèle et permettant de faire des opérations sur les fichiers stockés dans «Hadoop». Durant l'étape « **Map** », le nœud principal (master) prends les données d'entrées, les divises en problèmes plus petits, et les distribue aux nœuds de travail. Durant l'étape «**Reduce**» le nœud principal collecte les résultats de tous les sous-problèmes et les combine dans un résultat. Ce résultat est une réponse au problème posé initialement. (<a href="http://en.wikipedia.org/wiki/MapReduce">http://en.wikipedia.org/wiki/MapReduce</a>).

#### Hive

Hive est une infrastructure de «data warehouse» construite en top de « Hadoop ». Hive permet de faire de l'analyse, des résumés et des requêtes dans «Hadoop» (<a href="http://en.wikipedia.org/wiki/Apache\_Hive">http://en.wikipedia.org/wiki/Apache\_Hive</a>) avec un langage ressemblant de beaucoup à SQL. Ceci fait en sorte que c'est une solution très simple à implémenter.

Le problème avec Hive est qu'il est qu'il utilise les fonctions MapReduce qui sont très lentes à starter (aux alentour de 20 secs !!! selon nos test). De plus les jointures qui peuvent être faites sont limitées. De plus si jamais la structure des données de base change, les données devraient être réimportées d'oracle. Hive ne supporte pas (facilement) les opérations d'update ou d'effacement des données.

#### **Hbase**

HBase est une infrastructure par-dessus Hadoop, mais qui n'utilise pas MapReduce. HBase store les tables par colonnes (alors que traditionnellement ils sont stockées par rangées), ce qui luis permet de réduire grandement le temps de recherche (car on ne lits que les informations des colonnes dont on à besoin). HBase supporte les updates et les suppressions naturellement (avec son système de Versionnement). Par contre il n'y a aucun moyen simple et efficace de faire des jointures avec HBase.

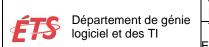
Voici la structure d'une entrée dans HBase :



Figure 3. HBase entrée. (http://www.larsgeorge.com/2009/10/hbase-architecture-101-storage.html)

Faire des recherches en connaissant la valeur d'une clé, où d'un préfixe de clé est très rapide. Tous les autres types de recherche vont demander de faire un scan complet de la table, qui est plus inefficace dans HBase que dans Hive.

Il serait très difficile d'utiliser cette solution directement pour le Dr.Eintracht, car ce dernier à besoin de rechercher en utilisant comme critère plusieurs colonnes sans ordre prédéfini. Par contre il pourrait être difficile si le problème du Dr. Eintracht puisse être réduit à un schéma en étoile (<a href="http://en.wikipedia.org/wiki/Star\_schema">http://en.wikipedia.org/wiki/Star\_schema</a>).



COURS LOG 792	DOCUMENT NO. 1.0	DATE 2012-09-24	version 1.0
TITRE		PAGE	PAGES
EndoMine - Rapport d'étape		11	38

#### Impala (Cloudera)

Le problème principal de Hive est le « *start up cost* » à cause de l'utilisation du MapReduce. Impala est un projet séparé développé par Cloudera. L'objectif est d'être un engin de recherche en temps réel pour «Hadoop»

(<a href="http://www.theregister.co.uk/2012/10/24/cloudera hadoop impala real time query/">http://www.theregister.co.uk/2012/10/24/cloudera hadoop impala real time query/</a>). En gros on peut le voir comme Hive, mais plus rapide (Important à noter il n'as pas toutes les possibilités de Hive mais dans notre cas d'utilisation on peut considérer comme tel. Voici le lien du papier de google sur lequel il est basé :

http://static.googleusercontent.com/external content/untrusted dlcp/research.google.com/en//pubs/archive/38125.pdf ).

Le projet est encore en stage beta. Pour cette raison nous allons éviter de l'utiliser à part si on n'a complètement pas le choix.

#### Pentaho:

Pentaho est un outil open source, pouvant intégrés divers sources de données ainsi que de faire des solutions de BI (Business Intelligence). Nous avons essayé les fonctions d'intégration est ils marchent plutôt bien. Nous n'avons pas encore eu le temps de tester les solutions BI (car nous ne savons pas si le problème de Dr. Eintracht peut être résumé à ces derniers). Le GUI pour les solutions BI n'est pas gratuit dans les majeures parties des cas.

# 4.1.2.3 Technologies Propriétaires

#### Datameer

À la conférence de Hadoop World, on a eu une démonstration d'un développeur de <u>DataMeer</u>. Le logiciel prend une approche de chiffrier et permet de faire des jointures, ainsi que des prévisualisations rapides (sur des sous-ensembles de données). Hadoop roule en back end.

- La version "Personal" (300\$/an, 1-node) devrait suffire pour le prototype.
- La version "Enterprise" permet d'utiliser un cluster Hadoop (approx. 9000\$/an).

#### Vertica, Greenplum, Teradata:

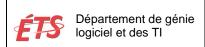
Des bases de données avec stockage par colonne. Toutes ces base de données sont chères, mais peuvent répondre au besoin de Mise à l'échelle.

### 4.1.3 Test Effectués sur différentes technologies

# 4.1.3.1 Installation de Hadoop et HBase :

Nous avons utilisé la distribution CDH4.1 de Cloudera car elle est stable et que l'installation est plus facile. De plus nous avons utilisé la version 1 de MapReduce (Mrv1), car la version 2 (Yarn) n'est pas encore considérée comme stable.

Voir **Annexe C** pour les instructions d'installation.



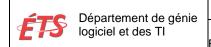
COURS LOG 792	DOCUMENT NO. 1.0	DATE 2012-09-24	VERSION 1 0
TITRE	1.0	PAGE	PAGES
EndoMine - Rappor	t d'étape	12	38

Le package Scoop est utilisé pour importer les données relationnelles dans HBase. Nous avons utilisé une base de données MySQL comme source de données initiale, mais ça devrait marcher avec la base de données Oracle.

#### 4.2 Recommandations

Étant donné les données disponibles, nous recommandons d'avoir les requêtes exactes de Dr. Eintracht. S'il est possible d'avoir une solution modélisé en tant que schéma en étoile ou autre modèle dimensionnel, nous allons utiliser Pentaho. Si ce n'est pas possible nous n'aurions pas d'autre choix que d'utiliser Hive avec les pertes de performance que ça encours.

Dans ce cas le cas d'utilisation de Hive les données des tables V\_P seront dé normalisés en une seule table, puis partitionnés par semaine. L'outil Visual SQL Query Builder pourrait être utilisé pour faire le Front End.



COURS	DOCUMENT NO.	DATE	VERSION
LOG 792	1.0	2012-09-24	1.0
TITRE		PAGE	PAGES
EndoMine - Rappor	t d'étape	13	38

# 5. LIVRABLES ET PLANIFICATION

# 5.1 Description des artéfacts

Nom de l'artefact	Description
Dlan de projet	Plan décrivant brièvement le projet, l'architecture / technologie et les
Plan de projet	personnes ressources, permettant d'initier le dialogue avec le client.
Échéancier	Plan décrivant la distribution des tâches dans le temps. Les tâches sont
(plan de travail)	sujettes à changement.
Proposition de projet	Document décrivant la proposition du projet, incluant la description, la
r roposition de projet	proposition, les risques et l'allocation du temps.
Schéma relationnel	Schéma relationnel des tables utiles de la base de données de
de production actuel	production. Pour des raisons de confidentialité, il est possible que
ас респисное посион	certaines tables soient renommées ou retirées du schéma.
	Document donnant plus de détails que la proposition du projet. Inclut la
Document de Vision	description et la portée. Inclut aussi le besoins du client, les
(B.R.S.)	caractéristiques du système ainsi que les demandes matérielles du
	système. Permet la discussion entre divers parties prenantes (« stake holders »).
Structure de	Définit la structure des fichiers dans Hadoop (HDFS, Hbase, Hive, etc.).
stockage des	Dominicia structure decinioro dano riddoop (ribr e, ribace, riive, etc.).
données	
Schéma architectural	Décrit les vues architecturales du système de « data mining ».
Rapport d'étape	Rapport intermédiaire pour le PFE.
Exigences systèmes	Exigences détaillées du système.
(S.R.S.)	
Scénarios de test /	Décrit les différents cas et scénarios de test pour confirmer la validité du
Cas de test	système pour le client.
Prototype de "data	Prototype fonctionnel utilisable par le client.
mining"	
Spécifications de	Spécifications pour les personnes désirant poursuivre le projet avec le
prochaine version	client.
Rapport final	Rapport final d'équipe à remettre à la fin du projet.
d'équipe	
Rapport final - DL	Rapport final de David Lauzon à remettre à la fin du projet.
Rapport final - AZ	Rapport final d'Anton Zakharov à remettre à la fin du projet.

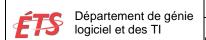
#### 5.2 Planification

Voir Annexe A.

# 6. RISQUES

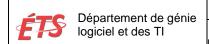
Les changements par rapport à la proposition de projet sont en gras.

D.		D 1 1 1114 /	Bertal at 1 at 1
Risque	Impact	Probabilité	Mitigation / atténuation
Misque	IIIIpaci	I I ODADIIILE	miligation / attenuation



COURS LOG 792	DOCUMENT NO. 1.0	DATE 2012-09-24	version 1.0
TITRE		PAGE	PAGES
EndoMine - Rappor	t d'étape	14	38

R1. La portée (« scope ») du projet est très ambitieuse.	Faible	Moyen	L'emphase du projet sera sur la documentation précise; ce qui facilitera les prochaines équipes qui travailleront sur le projet.  Le projet sera réalisé en plusieurs courtes itérations afin de pouvoir obtenir régulièrement de la rétroaction de la part du client.  Le prototype réalisé de répondra qu'à un seul besoin du client; ce qui fournira une preuve de concept et offrira un minimum de fonctionnalités au client.  Nous avons diminué le scope afin de répondre au besoin du Dr. Eintracht seulement.
R2. La confidentialité des données soit compromise.	Élevé	Faible	Afin de réduire le risque que les données confidentielles sur les patients soient divulgués, les informations personnelles sur les patients seront omises du système «Hadoop». Seul le Dr. Shaun aura accès aux données confidentielles.
R3. Les technologies utilisées sont mal connues.	Moyen	Faible	Lire la documentation et faire des exercices. De plus nous considérons la possibilité d'assister à la conférence «Hadoop World» en Octobre, ce qui inclue une activité de formation.  Nous avons une bonne idée des limitations des technologies.
R4. L'écosystème «Hadoop» ne répond pas aux besoins du client.	Moyen	Moyen	Nous allons commencer à s'informer sur «Hadoop» le plus tôt possible. De plus nous allons collaborer étroitement avec le client pour nous assurer qu'il comprend la solution et qu'il reçoit la formation nécessaire.  Nous nous somment informés. Les besoins du Dr. Shawn ne pourrons pas être résolus entièrement par aucun système sans ajout d'argent / machines supplémentaires.
R5. L'importation des données est inexacte.	Élevé	Faible	Nous allons développer un bon jeu de test, pour s'assurer que les données sont bien importées.  Nous allons probablement skipper cette partie.



COURS	DOCUMENT NO.	DATE	VERSION
LOG 792	1.0	2012-09-24	1.0
TITRE		PAGE	PAGES
EndoMine - Rappor	t d'étape	15	38

R6. Le système développé est moins performant que la solution Oracle existante.	Moyen	Élevé	Des discussions importantes avec le client seront réalisées afin de s'assurer que le système développé soit plus performant que le système existant.  Avec une seule machine c'est certains. Les systèmes distribués sont faits pour marcher avec plusieurs machines. Donc à court terme c'est inévitable.
R7. Les besoins du client sont mal compris.	Élevé	Faible	Nous avons planifié des rencontres hebdomadaires avec le client, en plus de réaliser plusieurs documents précisant le projet de manière itérative (Vision, SRS, Cas d'utilisation, etc.). Chaque document sera présenté au client afin de collecter une rétroaction permettant d'ajuster le projet aux besoins du client.  Les besoins sont biens compris mais ne sont pas réalisables tels quels. Par contre la performance pourrait être améliorée si Dr. Shawn nous présente les besoins les plus courants de façon plus spécifique.

# 7. RÉFÉRENCES CONSULTÉES

La liste des références ayant été consultées pour la réalisation du travail, la recherche de technologies existantes, et l'étude de faisabilité est placée à l'**Annexe B**.

# 8. TABLE DES MATIÈRES DU RAPPORT

Introduction

Problématique et contexte

Objectifs du projet

Méthodologie

Sommaire des travaux réalisés et recommandations

Sommaire des travaux réalisés

Définition des besoins du client

Exploration de différentes technologies

Test Effectués sur différentes technologies

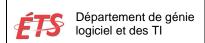
Choix technologique effectué.

Description de l'architecture.

Description d'installation des technologies.

Description du déploiement de la solution technologie chez le client.

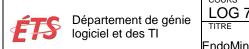
Référence pour modifier du schéma / importation des nouvelles données.



LOG 792	DOCUMENT NO. 1.0	DATE 2012-09-24	version 1.0
TITRE		PAGE	PAGES
EndoMine - Rappor	t d'étape	16	38

Recommandations
Livrables et planification
Description des artéfacts
Risques Restants
Références consultées
Bibliographie
Table des matières du rapport

Annexe A: Installation des technologies

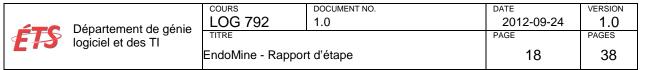


LOG 792	DOCUMENT NO. 1.0	DATE 2012-09-24	VERSION 1.0
TITRE		PAGE	PAGES
EndoMine - Rappor	t d'étape	17	38

# ANNEXE A : PLAN DE TRAVAIL RÉVISÉ

Indiquez les changements et ajouts en gras.

#	Commence	Termine	Efforts Estimés*	Tâches/Jalon	Livrable(s)/ Artéfacts	Responsabl e(s)
	13/09/2012	26/09/2012		Phase d'Inception - Itération I1		
1.1	14/09/2012	14/09/2012	1	Rencontre – au JGH avec le client		DL
1.2	14/09/2012	14/09/2012	1.5	Rencontre – professeur superviseur		DL
					Fiche de	
1.3	14/09/2012	20/09/2012	1	Définition de la fiche de renseignements	renseignements	DL
1.4	18/09/2012	28/09/2012	2	Planification du projet	Plan de projet	DL
2.1	20/09/2012	20/09/2012	4	Rencontre – au JGH avec le client		DL,AZ
				Identification et répartition des tâches du projet par		,
2.2	20/09/2012	24/09/2012	4	itération	Échéancier	DL
					Proposition de	
2.3	20/09/2012	24/09/2012	3	Définition de la proposition de projet	projet	DL,AZ
	27/09/2012	10/10/2012		Phase d'Inception - Itération I2		
3.1	27/09/2012	27/09/2012	3	Rencontre – au JGH avec le client		DL,AZ
					Schéma	
					relationnel de	
3.2	27/09/2012	03/10/2012	6	Analyse du système existant	production actuel	DL,AZ
				Recherche sur le forage de données à grande échelle (part		
3.3	27/09/2012	03/10/2012	5	1)		DL
4.4	04/10/2012	04/40/2042	1.5	Developed and CH average disease		DI 47
4.1	04/10/2012	04/10/2012	1.5	Rencontre – au JGH avec le client		DL,AZ
4.2	04/10/2012	10/10/2012	5	Identifier les problèmes du système existant		DL,AZ
	0.4/4.0/0.4.0	10/10/2010		Identifier les besoins et caractéristiques du nouveau		
4.3	04/10/2012	10/10/2012	8	système, et les prioriser		DL,AZ
	04/40/2042	40/40/2012	_	Recherche sur le forage de données à grande échelle (part		
4.4	04/10/2012	10/10/2012	5	2)	D	DL
4.5	01/10/2012	12/10/2012	5	Élicitation des besoins d'affaire.	Document de	AZ



					Vision (B.R.S.)	
	11/10/2012	21/10/2012		Phase d'Élaboration - Itération E1		
5.1	11/10/2012	11/10/2012	1.5	Rencontre – au JGH avec le client		DL,AZ
5.2	11/10/2012	17/10/2012	4	Concevoir les cas d'utilisations		DL,AZ
5.3	11/10/2012	17/10/2012	5	Identifier les exigences fonctionnelles, non-fonctionnelles, et contraintes de conception		AZ
5.5	11/10/2012	17/10/2012	J	Établir une stratégie de miroir de la BD de production		AL
5.4	11/10/2012	17/10/2012	5	actuelle		DL
5.5	11/10/2012	17/10/2012	5	Recherche sur le forage de données à grande échelle (part 3)		DL
6.1	18/10/2012	18/10/2012	1.5	Rencontre – au JGH avec le client		DL,AZ
6.2	18/10/2012	21/10/2012	5	Dénormaliser le schéma relationnel et sélection des champs importants		DL,AZ
6.3	18/10/2012	21/10/2012	5	Définir le schéma des relations Hadoop	Structure de stockage des données	DL,AZ
6.4	18/10/2012	21/10/2012	10	Concevoir une architecture pour le "data mining"	Schéma architectural	DL,AZ
6.5	18/10/2012	21/10/2012	5	Design des interface graphiques (recherche, résultats, etc.)	architectural	AZ
0.5	10/10/2012	21/10/2012		Recherche sur le forage de données à grande échelle (part		\
6.6	18/10/2012	21/10/2012	5	4)		DL
	22/10/2012	07/11/2012		Phase d'Élaboration - Itération E2		
7.1	22/10/2012	22/10/2012	1.5	Rencontre – au JGH avec le client		DL,AZ
7.2	22/10/2012	22/10/2012	1	Rencontre – professeur superviseur		DL,AZ
7.3	23/10/2012	25/10/2012	18	Formation/Conférence Hadoop World		DL,AZ
7.4	27/09/2012	29/10/2012	3	Réévaluation des objectifs du projet	Rapport d'étape	DL,AZ
7.5	15/10/2012	31/10/2012	5	Définition des exigences du client	Exigences systèmes (S.R.S.)	AZ



8.1	01/11/2012	01/11/2012	1.5	Rencontre – au JGH avec le client		DL,AZ
					Scénarios de test	
8.2	01/11/2012	09/11/2012	5	Élaboration des scénarios de test / cas de tests	/ Cas de test	DL,AZ
				Extraction/Conversion des données de BD production vers		
8.3	01/11/2012	07/11/2012	10	BD "miroir"		AZ
8.4	01/11/2012	07/11/2012	20	Installation de l'éco-système Hadoop	•	DL
	08/11/2012	21/11/2012		Phase de Construction - Itération C1		
9.1	08/11/2012	08/11/2012	1.5	Rencontre – au JGH avec le client		DL,AZ
				Implémentation de la structure Hadoop pour stocker les		
9.2	08/11/2012	14/11/2012	10	données		DL
9.3	08/11/2012	14/11/2012	5	Importer les données dans Hadoop		DL
9.4	08/11/2012	14/11/2012	5	Tester l'intégrité des données et des relations		DL
10.1	15/11/2012	15/11/2012	1.5	Rencontre – au JGH avec le client		DL,AZ
10.2	15/11/2012	21/11/2012	20	Implémentation de l'engin de recherche (part 1)		AZ
	22/11/2012	05/12/2012		Phase de Construction - Itération C2		
11.1	22/11/2012	22/11/2012	1.5	Rencontre – au JGH avec le client		DL,AZ
11.2	22/11/2012	22/11/2012	1	Rencontre – professeur superviseur		DL,AZ
11.3	22/11/2012	28/11/2012	5	Implémentation de l'engin de recherche (part 2)		AZ
11.4	22/11/2012	28/11/2012	5	Tester l'engin de recherche		DL
11.5	22/11/2012	28/11/2012	10	Implémentation de l'affichage et exportation des résultats		DL
11.6	22/11/2012	28/11/2012	5	Interface avec le système de visualisation des données		DL
11.7	22/11/2012	28/11/2012	15	Tester le prototype de "data mining"		AZ
12.1	29/11/2012	29/11/2012	1.5	Rencontre – au JGH avec le client		DL,AZ
	29/11/2012	30/11/2012	2	Démonstration d'un prototype de "data mining"	Prototype de "data mining"	DL,AZ
12.2	1 70/11/7017					

Département de génie logiciel et des TI		COURS LOG 792	DOCUMENT NO. 1.0	DATE 2012-09-24	VERSION 1.0
	TITRE EndoMine - Rapport d'étape		PAGE 20	PAGES 38	
		Lildolviille - Rappoi	i d clape	20	30

					d'utilisation	
	06/12/2012	14/12/2012		Phase de Transition - Itération T1		
13.1	06/12/2012	06/12/2012	1.5	Rencontre – au JGH avec le client		DL,AZ
					Spécifications de	
13.2	03/12/2012	14/12/2012	4	Documentation sur la continuation du projet	prochaine version	DL,AZ
13.3	06/12/2012	10/12/2012	2	Présentation oral (David Lauzon)		DL
13.4	06/12/2012	10/12/2012	2	Présentation oral (Anton Zakharov)		AZ
					Rapport final	
13.5	06/12/2012	12/12/2012	4	Rapport final d'équipe	d'équipe	DL,AZ
13.6	06/12/2012	12/12/2012	5	Rapport final individuel (David Lauzon)	Rapport final - DL	DL
13.7	06/12/2012	12/12/2012	5	Rapport final individuel (Anton Zakharov)	Rapport final - AZ	AZ
L	13/09/2012	14/12/2012		JALON 1 : Prototype de data mining		

Total d'heures David : 207,5 Total d'heures Anton : 185 Grand total: 392,5

<sup>\*</sup> En heures. \*\* En heures. Efforts réels.

	COURS	DOCUMENT NO.	DATE
∠ Département de génie	LOG 792	1.0	2012-09-24
logiciel et des TI	TITRE	PAGE	
	EndoMine - Rappor	t d'étape	21

# **ANNEXE B : RÉFÉRENCES**

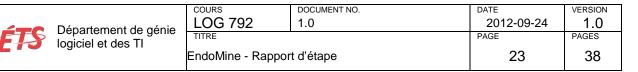
			Date	
	Date		de	
	de		consul	
Auteur	l'article	Titre de l'article	-tation	URL
			8	
			octobre	http://www.oracle.com/technetwork/java/javasebusiness/downloads/j
Oracle		Java SE 6 Downloads	2012	ava-archive-downloads-javase6-419409.html
	4		8	
Charles Tassfer	4 mai 2012	Install Organia Java IDI/ 7 in Librunty / Mint 40	octobre	http://www.inleans.com/dwwnol/pode/0
Charles Toepfer	2012	Install Oracle Java JDK 7 in Ubuntu / Mint 12	2012	http://www.iokom.com/drupal/node/9
			8 octobre	https://ccp.cloudera.com/display/CDH4DOC/Before+You+Install+CDH4+on+a+Single+Node#BeforeYouInstallCDH4onaSingleNode-
Cloudera		Before You Install CDH4 on a Single Node	2012	SupportedOperatingSystemsforCDH4
Cioddera		Defore Tod Iristali CDI 14 ori a Sirigle Node	8	Supported Operating Systems for ODI 14
			octobre	http://archive.cloudera.com/cdh4/cdh/4/hadoop/hadoop-project-
Cloudera		Deprecated Properties	2012	dist/hadoop-common/DeprecatedProperties.html
0.00.00.00			8	https://ccp.cloudera.com/display/DOC/CDH+Version+and+Packagin
			octobre	g+Information#CDHVersionandPackagingInformation-
Cloudera		CDH Version and Packaging Information	2012	CDHVersion4.1.0Packaging
		•	8	
			octobre	
Cloudera		CDH4 Installation	2012	https://ccp.cloudera.com/display/CDH4DOC/CDH4+Installation
			8	
		Installing CDH4 on a Single Linux Node in	octobre	https://ccp.cloudera.com/display/CDH4DOC/Installing+CDH4+on+a+
Cloudera		Pseudo-distributed Mode	2012	Single+Linux+Node+in+Pseudo-distributed+Mode
			8	
Olassalana		LIDana Installation	octobre	https://ccp.cloudera.com/display/CDH4DOC/HBase+Installation#HB
Cloudera		HBase Installation	2012	aseInstallation-InstallingHBase
			•	https://ccp.cloudera.com/display/CDH4DOC/ZooKeeper+Installation
Cloudera		ZooKeeper Installation	octobre 2012	#ZooKeeperInstallation-InstallingtheZooKeeperServerPackage
Oloudera	19	Zoortesper installation	8	#2001.coperinstaliation=installingthe2001.ceper0erverr ackage
Apache Sofware	février	HDFS File System Shell Guide (hadoop	octobre	
Foundation	2010	0.20.2)	2012	http://hadoop.apache.org/docs/r0.20.2/hdfs_shell.html
Apache Sofware		,	8	
Foundation		Apache HBase Book	_	http://archive.cloudera.com/cdh4/cdh/4/hbase/book.html

VERSION 1.0 PAGES

38



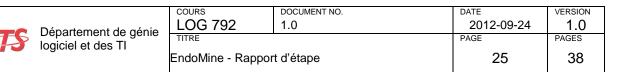
	Data		Date	
	Date de		de consul	
Auteur	l'article	Titre de l'article	-tation	URL
ratour	T un triolo	11110 00 1011010	2012	ONE.
	12			
	septem		8	
	bre	HBase Error – assignment of -ROOT- failure –	octobre	
Nathan	2011	Cant connect to web interface	2012	http://blog.nemccarthy.me/?p=110
	28		8	
	janvier	Quick install HBase in "pseudo distributed"	octobre	http://ria101.wordpress.com/2010/01/28/setup-hbase-in-pseudo-
Dominic Williams	2010	mode and connect from Java	2012	distributed-mode-and-connect-java-client/
			9	
Ala Davasa	9 avril	HBaseWD: Avoid RegionServer Hotspotting	octobre	http://blog.sematext.com/2012/04/09/hbasewd-avoid-regionserver-
Alex Baranau	2012	Despite Sequential Keys	2012	hotspotting-despite-writing-records-with-sequential-keys/
	novem	Hadaan Warld 2011, Advanged LIDaga	9	http://www.alidachara.pat/alaudara/badaan.warld 2011 advanced
Lars George	bre 2011	Hadoop World 2011: Advanced HBase Schema Design	octobre 2012	http://www.slideshare.net/cloudera/hadoop-world-2011-advanced-hbase-schema-design
Lais George	2011	Schema Design	9	Indase-scrienta-design
	june	Berlin Buzzwords June 2012: Advanced	octobre	http://www.slideshare.net/larsgeorge/hbase-advanced-schema-
Lars George	2012	HBase Schema Design	2012	design-berlin-buzzwords-june-2012
		Ŭ	11	,
			octobre	http://dev.mysql.com/doc/refman/5.5/en/connector-odbc-examples-
Oracle		Using Connector/ODBC with Microsoft Access	2011	tools-with-access.html
			11	
		Configuring a Connector/ODBC DSN on	octobre	http://dev.mysql.com/doc/refman/5.5/en/connector-odbc-
Oracle		Windows	2011	configuration-dsn-windows.html
			11	
Oragla		MyCOL, Download Commentar/ODBC	octobre	
Oracle		MySQL: Download Connector/ODBC	2011	http://dev.mysql.com/downloads/connector/odbc/
	4 mars	MySQL ODBC 32 vs 64 bit (answer from	octobre	
Justin Grégoire	2010	Justin Grégoire)	2011	http://stackoverflow.com/questions/2381906/mysql-odbc-32-vs-64-bit
Gustin Grogorio	2010	Jacan Gragona,	12	The particular of the wife of the particular of
			octobre	
Google		Google Refine	2012	http://code.google.com/p/google-refine/
			12	
			octobre	
OpenTSDB		Open Time Series Database	2012	http://opentsdb.net/



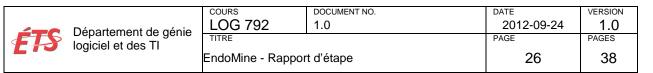
			Date	
	Date		de	
	de		consul	
Auteur	l'article	Titre de l'article	-tation	URL
			14	
	Août	No Relation: The Mixed Blessings of Non-	octobre	http://ianvarley.com/UT/MR/Varley_MastersReport_Full_2009-08-
Ian Varleys	2009	Relational Databases	2012	07.pdf
·	25		14	
	janvier	App Engine datastore tip: monotonically	octobre	http://ikaisays.com/2011/01/25/app-engine-datastore-tip-
Ikai Lan	2011	increasing values are bad	2012	monotonically-increasing-values-are-bad/
		HBASE-3551 Issue: Loaded hfile indexes		•
	20	occupy a good chunk of heap; look into	14	
	février	shrinking the amount used and/or evicting	octobre	
Michael Stack	2011	unused indices	2012	https://issues.apache.org/jira/browse/HBASE-3551
			14	http://search-
	24 mai	HBase, mail # user - a question	octobre	hadoop.com/m/hemBv1LiN4Q1/a+question+storefileIndexSize&subj
Gaojinchao	2011	storefileIndexSize	2012	=a+question+storefileIndexSize
-	24		15	http://search-
	mars		octobre	hadoop.com/m/nvbiBp2TDP/Stargate%252Bhbase&subj=Stargate+
Sreejith P. K.	2011	HBase, mail # user - Stargate+hbase	2012	hbase
	9		15	
	février		octobre	
Matteo Bertozzi	2011	HBase I/O: HFile	2012	http://th30z.blogspot.ca/2011/02/hbase-io-hfile.html?spref=tw
			15	
Apache Sofware			octobre	http://hbase.apache.org/xref/org/apache/hadoop/hbase/io/hfile/HFile.
Foundation		HFile Source Code	2012	html
	13		15	
	février		octobre	
Michael Stack	2009	HBASE-1200 Issue: Add bloomfilters	2012	https://issues.apache.org/jira/browse/HBASE-1200
	12		15	
	octobre		octobre	
Wikipedia	2012	Bloom filter	2012	http://en.wikipedia.org/wiki/Bloom_filter
			15	
	15 mai	StackOverflow: HBase MemStore and	octobre	http://stackoverflow.com/questions/10596717/hbase-memstore-and-
khan	2012	Garbage Collection	2012	garbage-collection
			15	
	5 août	StackOverflow: Where does HBase store all	octobre	http://stackoverflow.com/questions/6956400/where-does-hbase-
leon	2011	the row keys?	2012	store-all-the-row-keys



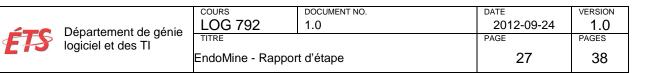
			Date	
	Date		de	
	de		consul	
Auteur	l'article	Titre de l'article	-tation	URL
	12		15	http://www.lauaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaa
Loro Coorgo	octobre	LIDaga Arabitaatura 101 Staraga	octobre	http://www.larsgeorge.com/2009/10/hbase-architecture-101-
Lars George	2009	HBase Architecture 101 - Storage	2012 15	storage.html
	janvier	StackOverflow: How to Scan HBase Rows	octobre	http://stackoverflow.com/questions/8961989/how-to-scan-hbase-
Panks	2012	efficiently	2012	rows-efficiently
1 diliko	2012	Cindiditaly	15	10W0 Ciriotetty
	3 août		octobre	http://permalink.gmane.org/gmane.comp.java.hadoop.hbase.user/28
AlexBaranau	2012	Re: How to query by rowKey-infix	2012	109
			15	
	mai		octobre	http://www.slideshare.net/cloudera/3-h-base-coprocessors-hbase-
Lars George	2012	HBaseCon: HBase Coprocessors	2012	con-may-2012
			16	
	22 mai		octobre	http://jimbojw.com/wiki/index.php?title=Understanding_Hbase_and_
Jimbojw	2008	Understanding HBase and BigTable	2012	BigTable
		LIVERTARIE VOLURAGE REREGRAANGE	16	
L lum amtalala		HYPERTABLE VS. HBASE PERFORMANCE EVALUATION II	octobre 2012	httm://b.va.autable.com/vvbv/ byva.autable.va.bb.co.2/
Hypertable		EVALUATION II		http://hypertable.com/why_hypertable/hypertable_vs_hbase_2/
	10		16	
	octobre	grokbase: [HBase-user] Hbase internally row	octobre	http://grokbase.com/t/hbase/user/10ab7vvfzy/hbase-internally-row-
William Kang	2010	location mechanism	2012	location-mechanism
	27		16	
Lars Hofhansl	janvier 2012	Scanning in HBase	octobre 2012	http://hadoop-hbase.blogspot.ca/2012/01/scanning-in-hbase.html
Lais Homansi	1		16	http://nadoop-ribase.biogspot.ca/2012/01/scarming-in-ribase.html
	février		octobre	
Quora	2011	How are bloom filters used in HBase?	2012	http://www.quora.com/How-are-bloom-filters-used-in-HBase
Quoid	2011	The Ware Steem Interested in Fibace.	16	map // www.querareem/riew are block micro about in ribace
		Culvert: A Robust Framework for Secondary	octobre	
Culvert	2012	Indexing	2012	https://github.com/booz-allen-hamilton/culvert
		-	16	
Apache Software		HBase API Docs: Package	octobre	http://hbase.apache.org/apidocs/org/apache/hadoop/hbase/coproces
Foundation		org.apache.hadoop.hbase.coprocessor	2012	sor/package-summary.html
	7		16	http://hadoop-hbase.blogspot.ca/2012/10/musings-on-secondary-
Lars Hofhansl	Octobr	Musings on Secondary Indexes	octobre	indexes.html



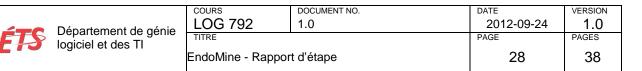
	_		Date	
	Date de		de consul	
Auteur	l'article	Titre de l'article	-tation	URL
Auteur	e 2012	Title de l'article	2012	OKE
	0 2012		2012	
			16	
	2 Juin		octobre	http://stackoverflow.com/questions/375194/how-to-design-hbase-
Yonatan	2011	how to design Hbase schema?	2012	schema
			16	
		Installing the Sqoop RPM or Debian	octobre	https://ccp.cloudera.com/display/CDH4DOC/Sqoop+Installation#Sqo
Cloudera		Packages	2012	opInstallation-installRPM
			16	
		Scoop User Guide: 7.2.11. Importing Data	octobre	http://archive.cloudera.com/cdh4/cdh/4/sqoop/SqoopUserGuide.html
Cloudera		Into HBase	2012	#_importing_data_into_hbase
Apache Software Foundation		Sqoop Developer's Guide v1.4.1-cdh4.1.0		http://archive.cloudera.com/cdh4/cdh/4/sqoop/SqoopDevGuide.html
Foundation		Squap Developer's Guide V1.4.1-cuil4.1.0	16	http://archive.cloudera.com/cdn4/cdn/4/sqoop/sqoopDevGuide.html
	2 mars		octobre	
J J Singh	2012	Sqoop installation tutorial	2012	http://jugnu-life.blogspot.ca/2012/03/sqoop-installation-tutorial.html
o o o migri	4			
	octobre	Google Groups: Cloudera Forum: Unable to		https://groups.google.com/a/cloudera.org/forum/?fromgroups=#!topic
Becky Benton	2011	load com.mysql.jdbc.Driver		/sqoop-user/pwdahVfAAAc
Booky Bornon		ioda commiyoqiijabolbrivor	16	70000 doon product in the to
			octobre	
Oracle		MySQL: Download Connector/J	2012	http://www.mysql.com/downloads/connector/j/
			16	
Apache Sofware		HBase Book : 16.1. Using existing ZooKeeper	octobre	
Foundation		ensemble	2012	http://hbase.apache.org/book/zookeeper.html
	40	Lilliana a communication Bata <b>7</b> and a second of 1941 Co	16	http://mail-archives.apache.org/mod_mbox/hbase-
N Kovavol	13 avril	Hbase-user mailing list: Zookeeper available	octobre 2012	user/201204.mbox/%3CCAPcDmSviyQXUG8u5dPmarjakFddndyCb
N Keywal	2012	but no active master location found	19	m6Pcgd8T6AR-aErw3g@mail.gmail.com%3E
	18 juin	Cloudera Developer Center: HBase Write	octobre	
Jimmy Xiang, Cloudera	2012	Path	2012	http://www.cloudera.com/blog/2012/06/hbase-write-path/
oning many, cloudera	16	1 400	19	The party of the party
	juillet	Configuring HBase Memstore: What You	octobre	http://blog.sematext.com/2012/07/16/hbase-memstore-what-you-
Alex Baranau	2012	Should Know	2012	should-know/



			Date	
	Date		de	
	de		consul	
Auteur	l'article	Titre de l'article	-tation	URL
	24			
	septem		19	
	bre	How partitioning, collecting and spilling work	octobre	
Alex Holmes	2012	in MapReduce	2012	http://grepalex.com/2012/09/24/map-partition-sort-spill/
				http://www.wired.com/wiredenterprise/2012/08/googles-mind-
			19	blowing-big-data-tool-grows-open-source-
	21 août	Google's Mind-Blowing Big-Data Tool Grows	octobre	twin/?utm_source=Contextly&utm_medium=RelatedLinks&utm_cam
Cade Metz (Wired)	2012	Open Source Twin	2012	paign=Previous
			19	
Dj Walker-Morgan (h-	21 août		octobre	http://www.h-online.com/open/news/item/Apache-to-Drill-for-big-
online)	2012	Apache to Drill for big data in Hadoop	2012	data-in-Hadoop-1671686.html
			19	
	9 août	Apache Drill: Interactive Analysis of Large-	octobre	http://wiki.apache.org/incubator/DrillProposal?action=AttachFile&do=
Tomer Shiran	2012	Scale Datasets	2012	view⌖=Drill+slides.pdf
			19	
Apache Software	9 août		octobre	
Foundation	2012	Apache Incubator Wiki : Drill Proposal	2012	http://wiki.apache.org/incubator/DrillProposal
			19	
			octobre	
Wikipedia		Pentaho	2012	http://en.wikipedia.org/wiki/Pentaho
			19	
			octobre	http://www.pentaho.com/resources/videos/67/pentaho-mapreduce-a-
Pentaho		Video Pentaho Presentation MapReduce	2012	major-league-baseball-use-case/
			19	
		Pentaho Community Edition (CE):	octobre	
Pentaho		Community Wiki Home	2012	http://wiki.pentaho.com/display/COM/Community+Wiki+Home
			19	
		Pentaho Community Edition (CE) : Latest	octobre	
Pentaho		Stable Builds	2012	http://wiki.pentaho.com/display/COM/Latest+Stable+Builds
Sergey Melnik, Andrey				
Gubarev, Jing Jing				
Long, Geoffrey Romer,				
Shiva Shivakumar, Matt			19	
Tolton, Theo Vassilakis		Dremel: Interactive Analysis of WebScale	octobre	http://static.googleusercontent.com/external_content/untrusted_dlcp/
Google, Inc.	2010	Datasets	2012	research.google.com/en//pubs/archive/36632.pdf
200gio, iiio.		24.400.0	120.2	1.000a.ogoogio.ooiii,oiii/pabo.aioiii/o/000o2.pai



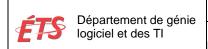
			Date	
	Date		de	
	de		consul	up.
Auteur	l'article	Titre de l'article	-tation	URL
	24 mai	Install Pentaho BI Server 4.5 on Ubuntu 12.04	19 octobre	http://akharahmad.com/2012/05/21/install_pantaha_hi_aanyar_1_5_an
Akbar Ahmed	24 mai 2012	LTS Desktop	2012	http://akbarahmed.com/2012/05/24/install-pentaho-bi-server-4-5-on-ubuntu-12-04-lts-desktop/
Akbai Alilleu	2012	LT3 Desktop	19	ubuntu-12-04-its-desktop/
	29 mai		octobre	http://akbarahmed.com/2012/05/29/install-kettle-4-3-0-on-ubuntu-12-
Akbar Ahmed	2012	Install Kettle 4.3.0 on Ubuntu 12.04 LTS	2012	04-lts/
7 HOAT 7 HITTOG	8	modal reduce more on obtained 1210 1 210	20.2	0.1.107
	décem		19	
	bre	Installing a Pentaho demo server on Ubuntu	octobre	
Sébastien Dejean	2006	6.10 Server Edition	2012	http://ubuntu-pentaho.blogspot.ca/
,	20			
	novem		19	
	bre	Pentaho Forums: Kettle repository - use and	octobre	http://forums.pentaho.com/showthread.php?65955-Kettle-repository-
codek	2008	how to create?	2012	use-and-how-to-create
			19	
	12 Sept		octobre	
bizcubed	2012	and Joins in Pentaho Data Integration	2012	http://www.youtube.com/watch?v=na6yRrhX5yo
			19	
<b>D</b>			octobre	http://wiki.pentaho.com/display/BAD/Configure+Pentaho+for+Cloude
Pentaho		Configure Pentaho for Cloudera CDH4	2012	ra+CDH4
			19	
Wikipadia		Classpath (Java)	octobre 2012	http://op.wikipodio.org/wiki/Closepoth / Joya)
Wikipedia	17	Classpatii (Java)	19	http://en.wikipedia.org/wiki/Classpath_(Java)
	février		octobre	
Mark Hall, John Paz	2012	HBase Input	2012	http://wiki.pentaho.com/display/EAI/HBase+Input
Wark Hall, John Faz	2012	Tibase input	21	mtp.//wiki.pentano.com/display/E/A//Tibase+mput
	août	Pentaho Hadoop Series: Big Data Analytics:	octobre	http://www.pentaho.com/resources/videos/25/hadoop-series-part-1-
James Dixon	2010	Part 1 - 5	2012	big-data-architecture/
			21	
		What are the advantages of Hadoop over	octobre	http://www.quora.com/What-are-the-advantages-of-Hadoop-over-
Quora		distributed RDBMS?	2012	distributed-RDBMS
			21	
			octobre	http://www-
IBM		Informix Features and Benefits	2012	01.ibm.com/software/data/informix/feature.html?S_CMP=rnav



	Data		Date	
	Date de		de consul	
Auteur	l'article	Titre de l'article	-tation	URL
7 10100	13		10.0.0	
	novem		21	
	bre		octobre	http://datawarehouse.ittoolbox.com/groups/strategy-planning/dw-
Glenn Engstrand	2011	OLAP Versus Big Data	2012	projectmanagement/olap-versus-big-data-4508739
			28	
		DataMeer: Enterprise: Analytics at the speed	octobre	
DataMeer		of business	2012	http://www.datameer.com/enterprise/index.html
	24		28	
Marcel Kornacker &	octobre	Cloudera Impala: Real-Time Queries in	octobre	http://blog.cloudera.com/blog/2012/10/cloudera-impala-real-time-
Justin Erickson	2012	Apache Hadoop, For Real	2012	queries-in-apache-hadoop-for-real/
			28	https://sep.elec.idexe.com/dienlec/IMDALA40DETADOC/Cleudexe.line
Cloudera		Cloudera Impala 1.0 Beta Documentation	octobre 2012	https://ccp.cloudera.com/display/IMPALA10BETADOC/Cloudera+Impala+1.0+Beta+Documentation
		·	2012	<del>                                     </del>
legrand_legrand	11 mai	SQLeo Visual Query Builder		http://sqleo.sourceforge.net/index.html
Tomer	2012	Announcing the MapR Hive ODBC Driver		http://www.mapr.com/blog/269?Itemid=78
Apache Software	2012	7 timounding the mapitalive obbo briver		mtp.//www.mapr.som/slog/2001.temid=70
Foundation		Hive JDBC Driver		https://cwiki.apache.org/Hive/hivejdbcinterface.html
	16		28	
	mars		octobre	
Rahul Patodi	2011	Hue Features	2012	http://www.technology-mania.com/2011/03/hue-features.html
	29			
	septem		28	
	bre	Using Different Reporting Frameworks with	octobre	http://wso2.org/library/articles/2012/09/using-different-reporting-
Sachini Jayasekara	2012	WSO2 Business Activity Monitor	2012	frameworks-wso2-business-activity-monitor
			28	
The Folince Foundation		New and Notable Features within BIRT 3.7	octobre 2012	http://www.colingo.org/hirt/phooniy/project/potable2.7.php
The Eclipse Foundation	27	I NEW AND INCIANCE FEATURES WITHIN DIKT 3.7	2012	http://www.eclipse.org/birt/phoenix/project/notable3.7.php
	septem		28	
	bre	phpHiveAdmin : Big data to Drive, Make	octobre	
xianglei	2012	easier for Hive	2012	http://www.phphiveadmin.net/
<u> </u>	26	-	28	
	janvier	StackOverflow: JavaScript Boolean Search	octobre	http://stackoverflow.com/questions/9022033/javascript-boolean-
yahelc	2012	Query Builder Interface Library?	2012	search-query-builder-interface-library



			Date	
	Date		de	
	de		consul	
Auteur	l'article	Titre de l'article	-tation	URL
			28	
		RedQueryBuilder - JavaScript SQL Query	octobre	
salk31		Builder UI	2012	http://redquerybuilder.appspot.com/
			28	
			octobre	http://www.developerextensions.com/index.php/extjs-grid-query-
Developer Extensions		Ext Grid Query Builder Example	2012	builder
			28	
	23 mai	Sencha:	octobre	http://www.sencha.com/forum/showthread.php?208444-
martinorth	2012	Ext.ux.window.VisualSQLQueryBuilder	2012	Ext.ux.window.VisualSQLQueryBuilder
			28	
			octobre	
martinorth		Visual SQL Query Builder	2012	http://www.cfsolutions.de/qb/
			28	
		Squel.js: Lightweight Javascript for building	octobre	
Ramesh Nair		SQL query strings	2012	http://hiddentao.github.com/squel/
			28	
			octobre	http://wiki.servoy.com/display/public/DOCS/Query+builder;jsessionid
Servoy		Servoy: Query Builder	2012	=7D0E08E4433B09432DC71E4F5584DE72



LOG 792	DOCUMENT NO. 1.0	DATE 2012-09-24	VERSION 1.0
TITRE		PAGE	PAGES
EndoMine - Rapport d'étape		30	38

### ANNEXE C: INSTALLATION DE HADOOP ET HBASE (DANS LINUX)

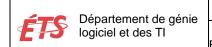
Oracle JAVA JDK Installation \_\_\_\_\_ cd Downloads/ chmod a+x jdk-6u34-linux-x64.bin sudo ./idk-6u34-linux-x64.bin sudo mv jdk1.6.0\_34//usr/lib/jvm/ Is -I /usr/lib/ivm sudo In -s /usr/lib/jvm/jdk1.6.0 34 /usr/lib/jvm/java-6-oracle sudo update-alternatives --install /usr/bin/java java /usr/lib/jvm/java-6-oracle/jre/bin/java 2 sudo update-alternatives --install /usr/bin/javac javac /usr/lib/jvm/java-6-oracle/bin/javac 1 sudo update-alternatives --install /usr/bin/javaws javaws /usr/lib/jvm/java-6-oracle/bin/javaws 1 sudo update-alternatives --config java sudo update-alternatives --install /usr/lib/mozilla/plugins/mozilla-javaplugin.so mozilla-javaplugin.so /usr/lib/jvm/java-6oracle/jre/lib/amd64/libnpjp2.so 1 sudo mkdir -p /opt/google/chrome/plugins sudo update-alternatives --install /opt/google/chrome/plugins/chrome-javaplugin.so chrome-javaplugin.so /usr/lib/jvm/java-6oracle/jre/lib/amd64/libnpjp2.so 1 sudo update-alternatives --config chrome-javaplugin.so sudo update-alternatives --config mozilla-javaplugin.so sudo vi /etc/profile.d/java.sh #!/bin/bash export JAVA\_HOME=\$( dirname \$( dirname \$( dirname \$( readlink -e /usr/bin/java ) ) ) ) export JDK HOME=\$JAVA HOME export JRE\_HOME=\$JAVA\_HOME/jre export PATH=\$JAVA HOME/bin:\$PATH sudo chmod +x /etc/profile.d/java.sh source /etc/profile.d/java.sh java -version

\_\_\_\_\_

Make sure the JAVA\_HOME environment variable is set for the root user

\_\_\_\_\_\_

sudo reboot



COURS	DOCUMENT NO.	DATE	VERSION
LOG 792	1.0	2012-09-24	1.0
TITRE		PAGE	PAGES
EndoMine - Rapport d'étape		31	38

[CTRL]+[ALT]+[F1] login as root env | grep JAVA\_HOME

\_\_\_\_\_\_

fix localhost problem

\_\_\_\_\_\_

# Remove entries relating to 127.0.1.1 < COMPUTER NAME > from /etc/hosts

# Make sure /etc/hostname matches the value in /etc/hosts

\_\_\_\_\_\_

#### **CDH4** Installation

\_\_\_\_\_\_

# Default data dir: /var/lib/hadoop-hdfs/cache/hdfs/dfs/data

sudo -u hdfs hdfs namenode -format

# Start HDFS services

sudo service hadoop-hdfs-namenode start

sudo service hadoop-hdfs-secondarynamenode start

sudo service hadoop-hdfs-datanode start

# Create Hadoop directory structure

sudo -u hdfs hadoop fs -mkdir /tmp

sudo -u hdfs hadoop fs -chmod -R 1777 /tmp

sudo -u hdfs hadoop fs -ls /

sudo -u hdfs hadoop fs -mkdir /var

sudo -u hdfs hadoop fs -mkdir /var/lib

sudo -u hdfs hadoop fs -mkdir /var/lib/hadoop-hdfs

sudo -u hdfs hadoop fs -mkdir /var/lib/hadoop-hdfs/cache

sudo -u hdfs hadoop fs -mkdir /var/lib/hadoop-hdfs/cache/mapred

sudo -u hdfs hadoop fs -mkdir /var/lib/hadoop-hdfs/cache/mapred/mapred

sudo -u hdfs hadoop fs -mkdir /var/lib/hadoop-hdfs/cache/mapred/mapred/staging

sudo -u hdfs hadoop fs -chmod 1777 /var/lib/hadoop-hdfs/cache/mapred/staging

sudo -u hdfs hadoop fs -chown -R mapred /var/lib/hadoop-hdfs/cache/mapred

sudo -u hdfs hadoop fs -ls -R /

# Start Map Reduce services

sudo service hadoop-0.20-mapreduce-jobtracker start

Département de génie logiciel et des TI
---

COURS LOG 792	DOCUMENT NO. 1.0	DATE 2012-09-24	VERSION 1.0
TITRE		PAGE	PAGES
EndoMine - Rapport d'étape		32	38

sudo service hadoop-0.20-mapreduce-tasktracker start

\_\_\_\_\_\_

#### Map Reduce Test

# Create a directory structure in HDFS for current user to run mapreduce jobs

sudo -u hdfs hadoop fs -mkdir /user/david

sudo -u hdfs hadoop fs -chown david /user/david

sudo -u hdfs hadoop fs -ls -R /

# Grep all config that starts with dfs from the hadoop config

hadoop fs -mkdir input

hadoop fs -put /etc/hadoop/conf/\*.xml input/

hadoop fs -ls input

hadoop jar /usr/lib/hadoop-0.20-mapreduce/hadoop-examples.jar grep input output 'dfs[a-z.]+'

hadoop fs -ls

hadoop fs -ls output

hadoop fs -cat output/part-00000 | head

\_\_\_\_\_

HBase Installation in Standalone mode

\_\_\_\_\_\_

sudo apt-get install hbase dpkg -L hbase

# Increase maximum number of open files for hdfs/hbase user sudo vi /etc/security/limits.conf

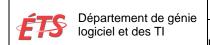
hdfs - nofile 32768 hbase - nofile 32768

sudo vi /etc/pam.d/common-session

session required pam limits.so

# Increase maximum number of files that can be served by a DataNode sudo cp /etc/hadoop/conf/hdfs-site.xml /etc/hadoop/conf/hdfs-site.xml.initial sudo vi /etc/hadoop/conf/hdfs-site.xml

operty>



LOG 792	DOCUMENT NO. 1.0	DATE 2012-09-24	VERSION 1.0
TITRE		PAGE	PAGES
EndoMine - Rapport d'étape		33	38

### 

# Restart HDFS services sudo service hadoop-hdfs-datanode stop sudo service hadoop-hdfs-secondarynamenode stop sudo service hadoop-hdfs-namenode stop sudo service hadoop-hdfs-namenode start sudo service hadoop-hdfs-secondarynamenode start sudo service hadoop-hdfs-datanode start

# Install HBase Master sudo apt-get install hbase-master sudo service hbase-master start

# Install HBase REST Interface sudo apt-get install hbase-rest

sudo cp /etc/hbase/conf/hbase-site.xml /etc/hbase/conf/hbase-site.xml.initial sudo vi /etc/hbase/conf/hbase-site.xml

<name>hbase.rest.port</name>
<value>60050</value>

sudo service hbase-rest restart

\_\_\_\_\_\_

Configuring HBase in Pseudo-distributed Mode

-----

Pseudo-distributed mode differs from standalone mode in that each of the component processes (e.g. HBase Master, Region Server and ZooKeeper peer) run in a separate JVM.

# Stop standalone HBase Master sudo service hbase-master stop

Département de gér logiciel et des TI	nie
---------------------------------------	-----

COURS	DOCUMENT NO.	DATE	VERSION
LOG 792	1.0	2012-09-24	1.0
TITRE		PAGE	PAGES
EndoMine - Rapport d'étape		34	38

# Set host to matches fs.default.name or fs.defaultFS in core-site.xml sudo vi /etc/hbase/conf/hbase-site.xml

# Create user for HBase sudo -u hdfs hadoop fs -mkdir /hbase sudo -u hdfs hadoop fs -chown hbase /hbase

# Installing the ZooKeeper Server Package and Starting ZooKeeper on a Single Server sudo apt-get install zookeeper-server sudo service zookeeper-server init sudo service zookeeper-server start

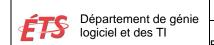
# Tell HBase to use separate JVM for ZooKeeper in # /etc/hbase/conf.dist/hbase-env.sh export HBASE\_MANAGES\_ZK=false

# Starting HBase Master sudo service hbase-master start

# Installing and Starting a HBase RegionServer sudo apt-get install hbase-regionserver sudo service hbase-regionserver start

# Verifying the Pseudo-Distributed Operation sudo jps sudo /usr/lib/jvm/jdk1.6.0\_34/bin/jps

> 32694 Jps 30674 HRegionServer 29496 HMaster



COURS	DOCUMENT NO.	DATE	VERSION
LOG 792	1.0	2012-09-24	1.0
TITRE		PAGE	PAGES
EndoMine - Rapport d'étape		35	38

### 28781 DataNode 28422 NameNode 30348 QuorumPeerMain

# Installing and Starting the HBase Thrift Server sudo apt-get install hbase-thrift

# Verifying HBase shell hbase shell

status 'detailed'

\_\_\_\_\_\_

#### **Testing HBase**

\_\_\_\_\_\_

\$ sudo -u hbase hbase shell HBase Shell; enter 'help<RETURN>' for list of supported commands. Type "exit<RETURN>" to leave the HBase Shell Version 0.92.1-cdh4.1.0, rUnknown, Sat Sep 29 11:55:59 PDT 2012

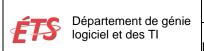
hbase(main):001:0> list TABLE 0 row(s) in 0.5590 seconds

hbase(main):002:0> create 'test', 'cf' 0 row(s) in 1.2210 seconds

hbase(main):004:0> list TABLE test 1 row(s) in 0.0400 seconds

hbase(main):006:0> put 'test', 'row1', 'cf:a', 'value1' 0 row(s) in 0.0900 seconds

hbase(main):007:0> put 'test', 'row2', 'cf:b', 'value2'
Auteurs : David Lauzon et Anton Zakharov



LOG 792	DOCUMENT NO. 1.0	DATE 2012-09-24	version 1.0
TITRE		PAGE	PAGES
EndoMine - Rapport d'étape		36	38

0 row(s) in 0.0310 seconds

hbase(main):008:0> put 'test', 'row3', 'cf:c', 'value3' 0 row(s) in 0.0280 seconds

hbase(main):009:0> scan 'test'

ROW COLUMN+CELL

row1 column=cf:a, timestamp=1349751879062, value=value1 row2 column=cf:b, timestamp=1349751884145, value=value2 row3 column=cf:c, timestamp=1349751892350, value=value3

3 row(s) in 0.0780 seconds

hbase(main):011:0> get 'test', 'row1'

COLUMN CELL

cf:a timestamp=1349751879062, value=value1

1 row(s) in 0.0150 seconds

hbase(main):012:0> disable 'test'

0 row(s) in 2.1320 seconds

hbase(main):013:0> list

TABLE test

1 row(s) in 0.0610 seconds

hbase(main):016:0> drop 'test' 0 row(s) in 1.4340 seconds

hbase(main):017:0> list

TABLE

0 row(s) in 0.0520 seconds

\_\_\_\_\_

Scoop Installation

\_\_\_\_\_

sudo apt-get install sqoop

Auteurs: David Lauzon et Anton Zakharov

	ÉTS	Département de génie logiciel et des TI	COURS LOG 792	DOCUMENT NO. 1.0
			TITRE EndoMine - Rappor	t d'étape

sqoop help sqoop version sqoop import # http://archive.cloudera.com/cdh4/cdh/4/sqoop/SqoopUserGuide.html # http://archive.cloudera.com/cdh4/cdh/4/sqoop/SqoopDevGuide.html

# Install MySQL JDBC Driver tar -xzf ~/Downloads/mysql-connector-java-5.1.22.tar.gz sudo cp mysql-connector-java-5.1.22/mysql-connector-java-5.1.22-bin.jar /usr/lib/sqoop/lib/chmod a+r /usr/lib/sqoop/lib/mysql-connector-java-5.1.22-bin.jar

# Should we set the connector path somewhere ??
#export SQOOP\_HOME="/home/hadoop/software/sqoop-1.3.0"
#export PATH=\$PATH:\$SQOOP\_HOME/bin

\_\_\_\_\_\_

### Scoop Test

\_\_\_\_\_\_

ERROR sqoop. Sqoop: Got exception running Sqoop: java.lang.RuntimeException: Could not load db driver class: com.mysql.jdbc.Driver

# Test JDBC connector and connection java -cp :/usr/lib/sqoop/lib/mysql-connector-java-5.1.22-bin.jar jdbctest 'jdbc:mysql://localhost/endomine access?user=root&password=...'

## # Initial test (hbase key != sql key)

sqoop import -libjars /usr/lib/sqoop/lib/mysql-connector-java-5.1.22-bin.jar --connect jdbc:mysql://localhost/endomine\_access -- username root -P --table vp\_lab\_order --split-by AA\_ID --hbase-table vp\_lab\_order --column-family d --hbase-row-key ID --hbase-create-table

DATE

PAGE

2012-09-24

37

VERSION

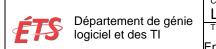
1.0

38

PAGES

list get 'vp\_lab\_order', "H3010002" get 'vp\_lab\_order', "H3010002", 'd:AA\_ID' disable 'vp\_lab\_order' drop 'vp\_lab\_order'

Auteurs: David Lauzon et Anton Zakharov



LOG 792 DOCUMENT NO.		DATE 2012-09-24	version 1.0
TITRE	PAGE	PAGES	
EndoMine - Rappor	t d'étape	38	38

# **ANNEXE D: DOCUMENT DE VISION**

(attaché à la page suivante)

Auteurs: David Lauzon et Anton Zakharov

# ENDOMINE

Projet de développement d'outils de forage de données de résultats de tests patients - endocrinologie, métabolisme et épidémiologie clinique

Projet #2012-076

Version : 0.9.1

Date d'émission : 4 octobre 2012 Date de révision : 11 octobre 2012

# Vision

Auteurs : Présenté à :



David Lauzon Anton Zakharov



Département d'endocrinologie

Document: Vision Date: 11 octobre 2012

# Table des matières

1	Introductio	
	1.1 Objecti	f
	1.2 Portée	
	1.3 Définiti	ons, acronymes et abréviations
	1.4 Référen	ces
	1.5 Langue	de rédaction
2	Positionner	nent
	2.1 Énoncé	du problème
	2.1.1	Problème 1
	2.1.2	Problème 2
	2.1.3	Problème 3
	2.2 Position	nnement du produit
3	Description	as des intervenants et des utilisateurs
	3.1 Résumé	$\epsilon$ des intervenants $(stakeholders)$
	3.2 Résumé	e des utilisateurs 🗋
	3.3 Environ	nement utilisateur
		aux besoins des intervenants et utilisateurs
	3.5 Alterna	tives et Compétition
4	Vue d'ensei	mble du produit
_		tive du produit
	_	é et Confidentialité
		aux avantages
		èses et dépendances
	HYP01	
	HYP02	
	HYP03	v .
	HYP04	
	HYP05	Conception, développement, et test
		s et installation
5	Consatániat	iques (features) du produit
J	FEA01	Configuration minimale pour ajouter d'autres machines
	FEA01	Code supportant le parallélisme
	FEA03	Distribution automatique des données pour supporter parallelisme
	FEA03	
		Supporter plusieurs requêtes en parallèles
	FEA05	Environnement de forage de données distinct des sources de données en
	· · · · · · · · · · · · · · · · · · ·	$\operatorname{production}$
	FEA06	Synchronisation des données automatique entre les sources de données et le
	\$	système EndoMine
		· ·

_[	Document: Vision	Date:	11 octobre 2012
	FEA09 Exportation des résultats de recherche	recherche	
6	Contraintes           VC01 Accessibilité            VC02 Confidentialité            VC03 Règles du JGH            VC04 Coût            VC05 Modification		
7	Gammes de qualité		1-
8	Attributs des caractéristiques		14
9	Autres exigences du produit9.1 Exigences du système9.2 Exigences de performance		
	DExigences de documentation  10.1 Manuel de l'utilisateur		
Li	iste des tableaux		
	Historique des Révisions Résumé des intervenants (stakeholders) Résumé des utilisateurs Besoins Avantages EndoMine Attributs des caractéristiques Légende : État des caractéristiques Légende : Bénéfice des caractéristiques Légende : Effort des caractéristiques Légende : Risque des caractéristiques Légende : Stabilité des caractéristiques Légende : Stabilité des caractéristiques Légende : Priorité des caractéristiques		
$\mathbf{T}_{i}$	able des figures		
	1 Modèle du domaine		

Version:

0.9.1

EndoMine (#2012-076)

Projet:

r rojec.	Endowine $(\#2012-070)$	A CIPIOII.	0.3.1	
Document:	Vision	Date:	11 octobre 2012	2
	nple d'une requête construite avec le Microsoft Access $Q$ siteweb : [2])		` _	18
	L 17			

Version: 0.9.1

EndoMine (#2012-076)

Projet:

Document: Vision Date: 11 octobre 2012

Table 1: Historique des Révisions

Date	Version	Description	Auteur
27 sept. 2012	v0.1	Sections 2.2, 2.3, 3.4, 3.7, 4.1	Anton Zakharov
27 sept. 2012	v0.2	Sections 1.3, 3.2, 3.3, 4.1, 4.3	David Lauzon
3 oct. 2012	v0.2.1	Formattage et mise en page	David Lauzon
3 oct. 2012	v0.3- preview	Section 1.2 et 3.7	Anton Zakharov
5 oct. 2012	v0.4	Section 4.2, 4.4, 4.5, 5-7	Anton Zakharov
7 oct. 2012	v0.5	Section 1.3-1.5, 3.2-3.4, 8-11	David Lauzon
8 oct. 2012	v0.5.1	Revue des sections à David	Anton Zakharov
9 oct. 2012	v0.9-rc1	Revue des sections à Anton	David Lauzon
11 oct. 2012	v0.9.1	Corrections suggérées par Fodil (2.1, 4.4, 4.5, 9.3).	David Lauzon

Document: Vision Date: 11 octobre 2012

### 1 Introduction

### 1.1 Objectif

Le but de ce document est de collecter, analyser, et définir les besoins et caractéristiques de haut niveau du système EndoMine. Il se concentre sur les fonctionalités recherchées par les parties prenantes, et explique pourquoi ces besoins existent. Les documents de cas d'utilisations (UC) et les spécifications des exigences logicielles (SRS) détaillent comment EndoMine satisfait ces besoins.

#### 1.2 Portée

Ce document de vision porte sur le développement et l'intégration du système de forage de données EndoMine. L'outil permettrait aux chercheurs du JGH de réaliser du forage de données à grande échelle sur la banque de données de tests biomédicaux. Les requêtes pourraient être faites à partir d'une interface utilisateur, tout en respectant les règles de sécurité et de confidentialité de l'hôpital.

La figure 1 présente les relations entre les principaux concepts dans le domaine du client. Une étoile signifie "plusieurs". Par exemple, la relation entre *Test Order* et *Test* se lit comme suis : "1 commande de tests peut comporter plusieurs tests", et la relation entre *Test* et *Test Result* se lit : "un même test n'a qu'un seul résultat". Les diagnostics (en jaune) et son lien avec le séjour n'existe pas encore dans le contexte du département de biochimie médicale, mais il est attendu qu'ils seront ajoutés en cours de projet. Se référer a la section 1.3 (Définitions, acronymes et abbréviations) pour l'explication de ces concepts et autres terminologie mentionné dans ce document.

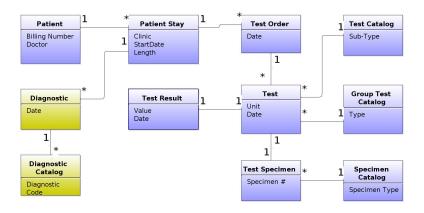


Figure 1 – Modèle du domaine

### 1.3 Définitions, acronymes et abréviations

BD, DB Base de données / banque de données / database
Diagnostic Associe le séjour d'un patient à un code de diagnostic.

Diagnostic Catalog Décris un diagnostic (ex : diabète). ÉTS École de Technologie Supérieure

Document: Vision Date: 11 octobre 2012

GÉLOG Software Engineering Research Laboratory

Group Test Catalog Décris le type général d'un test (ex : Testostérone).

JGH / HGJ Jewish General Hospital / Hôpital général juif

Patient Personne ayant un dossier (ou fiche) à l'hôpital JGH.

Patient Stay Séjour d'un patient à l'hôpital associé à une clinique en particulier

(ex : endocrinologie).

Specimen Catalog Décris la méthode du test (ex : urine, sang, etc.)

SQL Structured Query Language

Test Décris le test à effectuer sur un patient. Contient des informations

sur les unités du résultat de test (ex : g/mL).

Test Catalog Décris le type spécifique d'un test (ex : Testostérone-B,

Testostérone-C).

Test Order Commande d'un ou plusieurs tests effectuée à une date précise pour

un séjour en particulier.

Test Result Contient les informations relatives au résultat du test (ex : quantité

de mL de glucose).

Test Specimen Fait référence à l'échantillon de test (ex : l'éprouvette de sang).

### 1.4 Références

Voir le Plan de Projet pour la liste des artéfacts livrables de ce projet.

Les autres références citées se trouvent à la fin dans la section 10.2 intitulée Bibliographie.

### 1.5 Langue de rédaction

La langue de rédaction principale des documents sera le français. Toutefois, lorsque approprié, l'anglais pourrait être utilisé pour s'assurer la compréhension de l'ensemble des utilisateurs d'Endo-Mine.

### 2 Positionnement

### 2.1 Énoncé du problème

#### 2.1.1 Problème 1

Le problème de	l'extraction de statistiques complètes et précises sur l'utilisation des ressources médicales (scanner, microscope, tout ce qui ce trouve dans les laboratoires) est compliquée, voire impossible.
Cela affecte	les gestionnaires du laboratoire biomédical
dont l'impact est	la difficulté de prendre des décisions objectives d'achat et d'allocation de matériels.
Une bonne solution serait	d'ajouter les diagnostics médicaux et les relier aux tests médicaux dans le système de forage de données. Cela permettrait de prouver que des tests demandés par des médecins sont effectués inutilement, en vérifiant les tests avec les diagnostics.

Document: Vision Date: 11 octobre 2012

## 2.1.2 Problème 2

Le problème de	forer des données à grande échelle temporelle (données échelonnées sur plusieurs années) prend trop de temps et ralentis le système pour tous les utilisateurs.	
Cela affecte	les chercheurs et les utilisateurs d'équipements médicaux dans les laboratoires de biochimie	
dont l'impact est	que les chercheurs doivent soit : a) limiter la quantité de données re- cherchées ou b) attendre que les résultats d'une recherche complète soient disponibles.	
Une bonne solution serait	un système de forage de données rapide dont le traitement d'une re- cherche intensive soit quasi transparente pour les autres utilisateurs non concernés. De plus, la solution devrait s'adapter facilement à une quantité de données et à un nombre d'utilisateurs grandissants.	

### 2.1.3 Problème 3

Le problème de	la limite de 1GB de résultats de MS Access ne permet pas de re- chercher l'ensemble des données disponible (référence : [1]).	
Cela affecte	les chercheurs utilisant MS Access	
dont l'impact est	l'obligation de limiter la quantité de données traitées lors d'une même recherche.	
Une bonne solution serait	un système de forage de données permettant de faire des recherches sur l'ensemble des données de SoftLab, tout en conservant une fa- ciilité d'utilisation comme MS Access.	

# 2.2 Positionnement du produit

Pour	les chercheurs et gestionnaires des départements d'endocrinologie et diagnos- tique médicale de l'hopital général Juif	
qui	veulent une solution efficace pour faire des recherches dans leur banque de données.	
EndoMine	est une solution de forage de données	
qui	est efficace, rapide, et permet une mise à l'échelle (scaling) à faible coût.	
Contrairement à	une base données relationnelle traditionnelle, onéreuse, surchargée et lente,	
notre produit	est facile à utiliser, sauve du temps précieux dans la réalisation de recherches à grande échelle, tout en conservant un faible coût de possession (TCO).	

Document: Vision Date: 11 octobre 2012

# 3 Descriptions des intervenants et des utilisateurs

# 3.1 Résumé des intervenants (stakeholders)

Table 3 – Résumé des intervenants (stakeholders)

Nom	Description	Responsabilités
STK1. Dr. Elizabeth Mac	Chief Medical Biochemis-	Pilote (sponsor : states long term requi-
Namara	try (JGH)	rements).
STK2. Dr. Shaun Eintracht	Medical Biochemist,	Pilote (will operate and use the resul-
	Dept. of Diagnostic	ting data mining system).
	Medecine (JGH)	
STK3. Dr. Mark Trifiro	Chief Endocrinology and	Client (sponsor : states long term requi-
	Metabolism (JGH)	m rements).
STK4. Dr. Sami Suissa	Chief Clinical Epidemio-	Client (sponsor : states long term requi-
	logy (JGH)	m rements).
STK5. Dr. Alain April	Directeur du GÉLOG	Provide Information Technology solu-
	(ÉTS)	tions and guide ÉTS students.
STK6. Chris Polykandriotis	IT Specialist (JGH)	Help with IT issues.

## 3.2 Résumé des utilisateurs

Table 4 – Résumé des utilisateurs

Nom USR1. Utilisateur privilégié	Description  Utilisateur ayant accès à la base de données de production Oracle, et au système Endo-Mine. Il s'agit des utilisateurs du département de biochimie du JGH.	Responsabilités  - Faire des recherches spécifiques.  - Produire des rapports de recherche  - Gérer l'utilisation de ressources médicales.	Intervenant STK1 (Dr. Mac Namara), STK2 (Dr. Eintracht)
USR2. Utilisateur restreint	Utilisateur externe au département de biochimie ayant accès seulement au système Endo-Mine (donc ils n'ont pas accès aux informations confidentielles des patients). Il s'agit des utilisateurs des autres départements du JGH.	– Faire des recherches spécifiques.	STK3 (Dr. Trifiro), STK4 (Dr. Suissa)

Document: Vision Date: 11 octobre 2012

### 3.3 Environnement utilisateur

- Poste de travail ayant un accès intranet au serveur d'EndoMine.

## 3.4 Principaux besoins des intervenants et utilisateurs

Table 5: Besoins

Besoin	Priorité	Préoccupations	Solution actuelle	Solution proposée
N1. Confi-	Critique	Informations confiden-	Mesures de sécu-	Voir la section 4.2 ('Sé-
dentialité des		tielles protégées contre	rité inconnues.	curité et Confidentiali-
données		l'accès de personnes non		té').
		autorisées.		
N2. Forage de	Critique	Comme le forage est	Aucune	Séparation distincte des
données n'in-		traité directement sur		environnements de pro-
terférant pas		l'environnement de		duction et de forage de
avec la collec-		collection des données,		données.
tion de don-		des pannes surviennent		
nées		bloquant temporaire-		
		ment l'accès à tous les		
		utilisateurs.		
N3. Outil	Critique	Courte formation pour	Dr. Eintracht	Interface similaire à MS
simple pour		apprendre l'outil. Au-	utilise MS Access	Access Query Builder.
créer des		cune connaissance de	Query Builder	(voir la Figure 3 de
requêtes		SQL requise.		1'10.2)
N4. Mise à	Important	Accommoder un nombre	Incapable de ré-	Extensibilité horizontale
l'échelle facile		de données et d'utilisa-	pondre à ce be-	graduelle en ajoutant
du système		teurs grandissant, sans	soin.	d'autres ordinateurs au
		que les coûts ne de-		système. La configura-
		viennent exponentiels.		tion à modifier serait mi-
				nimale.
N5. Forage	Important	Pouvoir rechercher	La mémoire de	L'architecture d'En-
de données		toutes les données	MS Access li-	doMine n'as pas cette
sans limite de		disponibles.	mite la taille des	restriction.
taille			recherches. [1]	

## 3.5 Alternatives et Compétition

A notre connaissance il n'existe pas de solution clé en main pour faire directement du forage de données sur la structure de la base de données actuelle. Une implémentation spécifique aux besoins du JGH est donc nécessaire.

Document: Vision Date: 11 octobre 2012

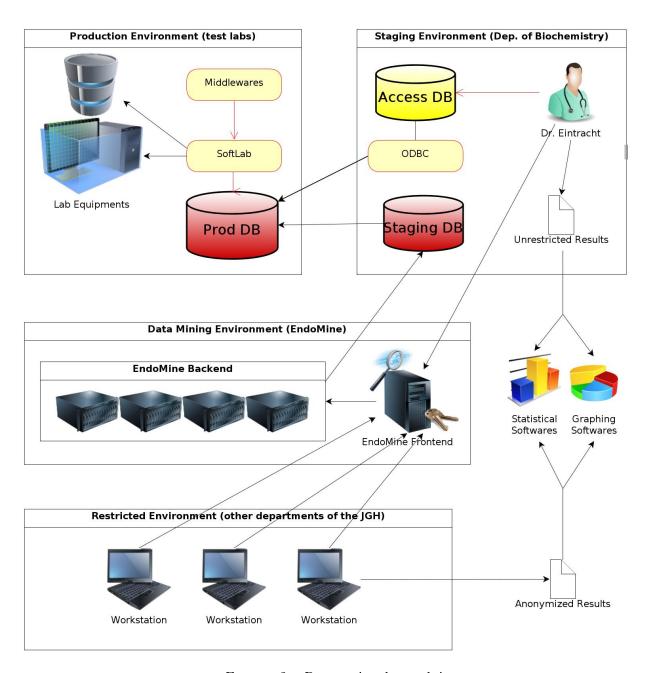


Figure 2 – Perspective du produit

Document: Vision Date: 11 octobre 2012

# 4 Vue d'ensemble du produit

### 4.1 Perspective du produit

La figure 2 situe EndoMine en perspective avec les autres produits mis en place dans l'environnement du JGH.

- Le Production Environment représente l'environnement où les données sont collectées et la BD
   Oracle de production est populée par le logiciel SoftLab.
- La Staging DB est une copie miroir de la BD de production dont les informations confidentielles ont été anonymisées. Seul le département de biochimie a accès au Staging Environment.
- Les utilisateurs effectueront le forage de données sur le système EndoMine, et pourront par la suite importer leur résultats dans des logiciels de statistiques et de graphiques.

### 4.2 Sécurité et Confidentialité

Cette section décrit les différentes mesures à prendre pour protéger la confidentialité des données :

- L'accès au système EndoMine serait restreint selon les protocoles de sécurité du JGH.
- Les informations confidentielles (noms des patients, numéro de RAMQ, numéro d'hôpital, etc.) seront absentes du système EndoMine.
- EndoMine contiendra la clef anonymisée des patients, mais la correspondance entre la clef anonymisée et la fiche d'hôpital du patient ne sera pas accessible par le système EndoMine.
- Voir la section 4.4 ('Hypothèses et dépendances') pour plus d'informations.

### 4.3 Principaux avantages

Table 6: Avantages EndoMine

Bénéfices pour le client	Caractéristiques		
	${\it correspondantes}$		
N1. Confidentialité et sécurité des données	HYP01, HYP03,		
	HYP04, FEA05,		
	FEA10, VC01,		
	VC02, VC03		
N2. Forage de données n'interférant pas avec la collection de données	HYP02, FEA05,		
	FEA06, VC05		
N3. Outil simple pour créer des requêtes	FEA07, FEA08,		
	FEA09, FEA11,		
	FEA12, VC04		
N4. Mise à l'échelle facile du système	FEA01, FEA02,		
	FEA03, FEA04,		
	VC04		
N5. Forage de données sans limite de mémoire	FEA02, FEA03,		
	FEA11		

Document: Vision Date: 11 octobre 2012

### 4.4 Hypothèses et dépendances

### HYP01 BD de staging

La BD de staging sera fournie par le JGH, et sera déja anonymisée (voir la section 4.2).

### **HYP02** Synchronisation

La BD de staging sera synchronisée automatiquement une fois par jour (ou toute autre fréquence jugée acceptable par le client).

### HYP03 Anonymisation

Les scripts pour anonymiser / désanonymiser seront mis à la disposition du Dr. Eintracht par le JGH.

### HYP04 Prototype

Lors du prototype, seul le Dr. Eintracht aura accès au système EndoMine installé au JGH.

### HYP05 Conception, développement, et test

Les ressources de l'ÉTS doivent avoir accès aux données anonymisées et au système pilote pour les besoins de conception, développement et test.

#### 4.5 Licences et installation

- 1. EndoMine sera implémenté avec logiciels libres ayant prouvé leur efficacité sur des sites tels que Facebook, Amazon et Yahoo. Donc aucune license à acheter.
- 2. Le code développé pour le logiciel EndoMine sera la propriété commune du JGH, de l'ÉTS et des développeurs ayant travaillé sur EndoMine. Les propriétaires du logiciel pourront installer une copie du logiciel sur d'autres serveurs et y effectuer des modifications.
- 3. Les données provenant de la base de données de production Oracle, ainsi que les résultats de recherches réalisés avec EndoMine sont la propriété exclusive du JGH.
- 4. Le système EndoMine sera installé sur le serveur du JGH par les développeurs de l'ÉTS.
- 5. Aucune installation requise sur les postes de travails des chercheurs.

# 5 Caractéristiques (features) du produit

### FEA01 Configuration minimale pour ajouter d'autres machines.

L'objectif est de monter le système de telle façon que le coût de configuration / stabilité avec l'ajout d'une nouvelle machine parallèle soit minimal.

#### FEA02 Code supportant le parallélisme.

Document: Vision Date: 11 octobre 2012

En lien avec FEA01, le forage de données devrait pouvoir exploiter le parallélisme du système de machines distribués.

### FEA03 Distribution automatique des données pour supporter parallelisme

Une fois une nouvelle machine installée, le système devrait balancer automatiquement (le plus possible) la charge de traitement sur la nouvelle machine.

### FEA04 Supporter plusieurs requêtes en parallèles

Plusieurs requêtes d'un même (ou plusieurs) utilisateur(s) devraient pouvoir être exécutés en parallèle.

# FEA05 Environnement de forage de données distinct des sources de données en production

L'utilisation du système développe ne devrait avoir aucun impact sur le contenu ou le fonctionnement des sources de données en production.

# FEA06 Synchronisation des données automatique entre les sources de données et le système EndoMine

EndoMine devrait pouvoir se synchroniser automatiquement avec les sources de données (dans l'environnement de *staging*) à intervalles régulières configurables.

#### FEA07 Générateur de requêtes intégré

Un générateur de requêtes permettant de faire des requêtes d'une manière interactive et simple devrait être fournis. EndoMine s'inspirera de *Microsoft Access Query Builder* (voir la Figure 3 de l'10.2).

### FEA08 Recherche par filtrage incrémentiel de la requête originale

Modification d'une requête en ajoutant des conditions de recherche, qui ne faisait pas partie de la requête originale. C'est-à-dire de pouvoir filtrer les résultats d'une recherche avec des nouveaux critères, sans que le système aille besoin de ré-exécuter la requête originale. Effectuer un filtrage secondaire devrait prendre une fraction du temps de la requête originale. Par exemple, si on recherche les résultats de Glucose en 2010, on devrait pouvoir ajouter les résultats de Glucose en 2009 et/ou ajouter les résultats de Fructose.

### FEA09 Exportation des résultats de recherche

Les résultats de recherche devraient pouvoir être exportés selon au moins 1 format d'échange de fichier d'un logiciel d'analyse statistique.

### FEA10 Exportation des clés anonymisées

Document: Vision Date: 11 octobre 2012

Permettre de facilement exporter seulement l'ensemble unique des clés anonymisées des patients présents dans les résultats d'une recherche.

### FEA11 Ajout de champs supplémentaire à des résultats de recherche

À partir de résultats de recherche, il doit être possible d'ajouter des champs qui n'étaient pas inclus dans la sélection. Par exemple, on pourrait vouloir ajouter le champs sexe et âge du patient. Tous les champs disponibles dans EndoMine devraient pouvoir être ajoutés de cette façon.

#### FEA12 Trier les résultats de recherche

Pouvoir trier les résultats de recherche selon n'importe quel champ inclus dans la recherche.

### 6 Contraintes

#### VC01 Accessibilité

Le système EndoMine ne doit pas être accessible à l'extérieur du réseau du JGH.

#### VC02 Confidentialité

Les utilisateurs restreints ne peuvent accéder aux informations confidentielles sur les patients. Voir la section 4.2 (Sécurité et Confidentialité).

### VC03 Règles du JGH

Les règles de sécurité informatique et de confidentialité du JGH doivent être respectées.

### VC04 Coût

La réalisation et l'installation du prototype doit être de faible coût.

### VC05 Modification

La base de données de production ne doit pas être modifiée par EndoMine.

## 7 Gammes de qualité

Le forage de données doit être au minimum deux fois plus rapide que le système actuel.

# 8 Attributs des caractéristiques

Le tableau suivant permet :

Document: Vision Date: 11 octobre 2012

- Au client de prendre connaissance des efforts et risques associés au projet.

- Aux développeurs de jauger les bénéfices et priorité que le client associe à chaque caractéristique.
- Et d'identifier les caractéristiques les plus susceptibles de changer dans le futur.

La légende des valeurs possibles pour chacune des colonnes est présenté à l'10.2.

Table 7: Attributs des caractéristiques

Caractéristiques	État	Bénéfice	Effort	Risque	Stabilité	Priorité
FEA01. Configuration minimale pour	Proposé	Moyen	Moyen	Moyen	Élevé	Important
ajouter d'autres machines						
FEA02. Code supportant le parallélisme	Proposé	Moyen	Faible	Faible	Élevé	Important
FEA03. Distribution automatique des	Proposé	Moyen	Faible	Faible	Élevé	Important
données pour supporter parallelisme						
FEA04. Supporter plusieurs requêtes en	Proposé	Élevé	Faible	Faible	Élevé	Important
parallèles						
FEA05. Environnement de forage de	Proposé	Élevé	Moyen	Faible	Moyen	Critique
données distinct des sources de données						
en production						
FEA06. Synchronisation des données	Proposé	Élevé	Moyen	Moyen	Faible	Important
automatique entre les sources de don-						
nées et le système EndoMine						
FEA07. Générateur de requêtes intégré	Proposé	Élevé	Élevé	Élevé	Faible	Critique
FEA08. Recherche par filtrage incré-	Proposé	Moyen	Élevé	Élevé	Faible	Critique
mentiel de la requête originale						
FEA09. Exportation des résultats de re-	Proposé	Élevé	Moyen	Faible	Faible	Critique
cherche						
FEA10. Exportation des clés anonymi-	Proposé	Moyen	Faible	Faible	Moyen	Important
sées						
FEA11. Ajout de champs supplémen-	Proposé	Moyen	Élevé	Élevé	Faible	Critique
taire à des résultats de recherche						
FEA12. Trier les résultats de recherche	Proposé	Moyen	Moyen	Faible	Moyen	Utile

# 9 Autres exigences du produit

### 9.1 Exigences du système

ES1. Réquis minimaux de l'ordinateur exécutant la recherche :

- Fureteur Firefox 8 ou équivalent.
- 2 GB de mémoire vive (RAM).
- Suite de bureautique Office installée (optionnel).
- Logiciel(s) d'analyse(s) statistique(s) (optionnel).

Document: Vision Date: 11 octobre 2012

## 9.2 Exigences de performance

??? SVP. Précisez si cette section s'applique.

# 10 Exigences de documentation

### 10.1 Manuel de l'utilisateur

EMU1. Un manuel d'utilisateur simple et complet, permettra à l'utilisateur d'apprendre rapidement le fonctionnement de l'interface de création de requêtes.

## 10.2 Guides d'installation, de configuration, et fichier à lire

EGCF1. Documenter les configurations requises pour installer et configurer une machine supplémentaire pour le système EndoMine.

EGCF2. Documenter les configurations requises pour ajouter une nouvelle source de données ou table supplémentaire (si le temps le permet).

Document: Vision Date: 11 octobre 2012

# 11 Bibliographie

[1] Access 2010 specifications. Microsoft. Consulté le 3 octobre 2012, à http://office.microsoft.com/en-us/access-help/access-2010-specifications-HA010341462.aspx

[2] ArchSummary Query. 10 fév. 2010. Jessica Iannone. Consulté le 3 octobre 2012, à http://jessica-iannone.com/img/interactive/ArchSummQry.jpg

Document: Vision Date: 11 octobre 2012

# Microsoft Access Query Builder

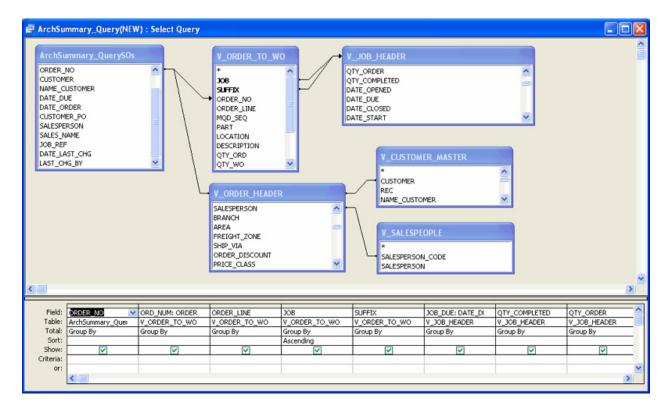


Figure 3 – Exemple d'une requête construite avec le Microsoft Access Query Builder (image tirée d'un siteweb : [2])

Document: Vision Date: 11 octobre 2012

# Attributs des caractéristiques

	Table 8 – Légende : État des caractéristiques			
Proposé	La caractéristique est proposée, mais n'a pas encore été approuvée par les parties prenantes.			
Approuvé	La caractéristique est approuvée par les parties prenantes.			
Incorporé	La caractéristique est incluse dans le produit.			
	Table 9 – Légende : Bénéfice des caractéristiques			
Faible	La caractéristique apporte peu de valeur ajoutée au produit et n'est pas nécessaire à son bon fonctionnement.			
Moyen	La caractéristique apporte une valeur ajoutée additionnelle au produit, mais n'est pas critique à son bon fonctionnement.			
Élevé	La caractéristique apporte une valeur ajoutée importante au produit et est essentielle à son bon fonctionnement ou à la réalisation de ses tâches.			
	Table 10 – Légende : Effort des caractéristiques			
Faible	La réalisation de la caractéristique nécessite un effort de moins de 20 heures- personnes.			
Moyen	La réalisation de la caractéristique nécessite un effort entre 20 et 40 heures- personnes.			
Élevé	La réalisation de la caractéristique nécessite un effort de plus de 40 heures- personnes.			

Document: Vision Date: 11 octobre 2012

	Table 11 – Légende : Risque des caractéristiques			
Faible	La technologie utilisée et la méthode d'implémentation sont connues et bien maîtrisées.			
Moyen	La technologie utilisée est récente ou la méthode d'implémentation nécessite une attention particulière.			
Élevé	La technologie utilisée est nouvelle et peu éprouvée ou la méthode d'implémentation est complexe et demande une analyse plus complète.			
	Table 12 – Légende : Stabilité des caractéristiques			
Faible	Les exigences concernant la caractéristique ont de fortes chances de changer ou le bon fonctionnement de la caractéristique a un impact critique sur le fonctionnement général du système et peut compromettre son exécution.			
Moyen	Les exigences concernant la caractéristique sont susceptibles de changer ou le bon fonctionnement de la caractéristique a un impact sur le fonctionnement général du système sans toutefois compromettre son exécution.			
Élevé	Les exigences concernant la caractéristique ont peu de chance de changer et le bon fonctionnement de la caractéristique n'a pas d'impact sur le fonctionnement général du système.			
	Table 14 – Légende : Priorité des caractéristiques			
Utile	La caractéristique apporte des fonctionnalités accessoires au système. Son inclusion dans le produit a peu d'impact sur la satisfaction du client et sur l'utilisation du système.			
Important	La caractéristique apporte des fonctionnalités supplémentaires au système. Son inclusion dans le produit peut influencer la satisfaction du client, mais son absence n'empêche pas l'utilisation du système.			
Critique	La caractéristique est primordiale au fonctionnement du système. Il est nécessaire de l'inclure en priorité dans le produit pour assurer la totale satisfaction du client et son absence pourrait empêcher l'utilisation du système.			