

Fouillez facilement dans votre système Big Data

Olivier TAVARD

Introduction

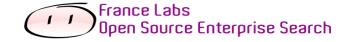
A propos de moi :

- Cofondateur de la société France Labs
- Développeur (principalement Java)
- Formateur en technologies de moteurs de recherche (Solr, Elasticsearch)

A propos de France Labs :

- Startup créée en 2011
- Experts en technos de search (consulting)
- Partenaire officiel de LucidWorks
- Editeur de la solution Datafari





Plan

- Rappels sur Hadoop
- Présentation (rapide) de Solr
- Intégration d'Hadoop et de Solr
- Démo



Hadoop

- Créé par Doug Cutting en 2004
- A l'origine pour Lucene : Projet Nutch
- Framework open source
- Inspiré par les papiers sur Google Map Reduce et Google File System









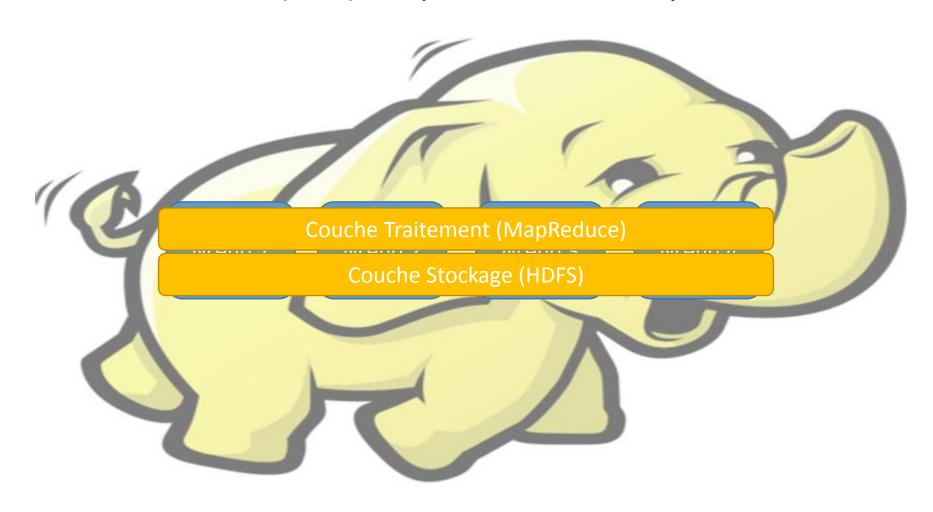
Hadoop

- Traiter des grands volumes de données en un minimum de temps
- Stocker des immenses volumes de données
- Fonctionne sur machines de configuration faible et peu coûteuses



Hadoop

Architecture (très) simplifiée de Hadoop v1



Hadoop

Couche stockage (HDFS)

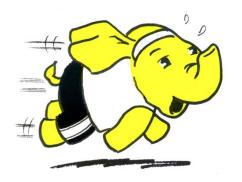
- Système de fichier distribué
- Utilise les disques « normaux » de chacun des nœuds...
 - ...pour donner l'impression de n'utiliser qu'un seul énorme disque
- Fiable (données répliquées 3 fois)



Hadoop

Couche traitement (MapReduce)

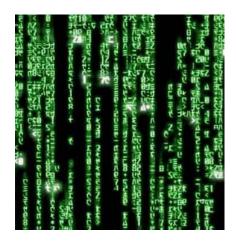
- Hadoop exécute des jobs Map Reduce
- Un algorithme doit donc implémenter:
 - Un mapper
 - Un reducer
- Pour être executé de manière distribuée par Hadoop en tant que job



Hadoop

Mapper

- Prend des données en entrée
- Et les transforme en paire de clé valeur





Hadoop

Reducer

Combine les valeurs ayant la même clé

Clé 1 : Val2 Clé 2 : Val2

Clé 1 : Val3

Clé 2: Val4



Hadoop

Un exemple simple?

- On a en entrée plusieurs textes
- On veut compter le nombre d'occurrences des mots
- De manière distribuée
- En utilisant un job Map Reduce

Solr

Apache Solr est une couche web open source basée sur la librairie de recherche Apache Lucene. Elle ajoute des fonctionnalités de serveur web, du faceting.

Datafari

Datafari est notre solution open source de recherche d'entreprises, basée sur Lucene, Solr, et le framework de connecteurs Apache ManifoldCF.

Lucene

Apache Lucene est un moteur de recherche et d'indexation, en open source.



Hadoop

Solr

Apache Solr est une couche web open source basée sur la librairie de recherche Apache Lucene. Elle ajoute des fonctionnalités de serveur web, du faceting.



Apache: 1

Apache: 1

Solr:1

Lucene: 1

...



Datafari est notre solution open source de recherche d'entreprises, basée sur Lucene, Solr, et le framework de connecteurs Apache ManifoldCF.

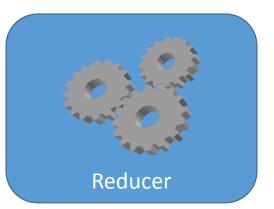


Apache: 1

Solr:1

Lucene: 1

•••



Apache: 4

Solr: 2

Lucene: 3

•••



Apache Lucene est un moteur de recherche et d'indexation, en open source.



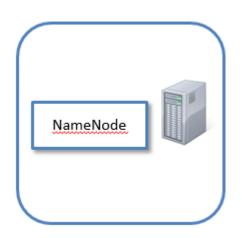
Apache : 1 Lucene : 1

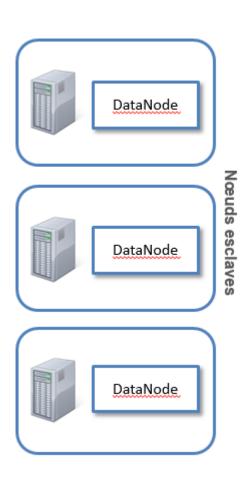
•••

Hadoop 2

HDFS 2:

- Haute disponibilité NameNode
- Support de NFS



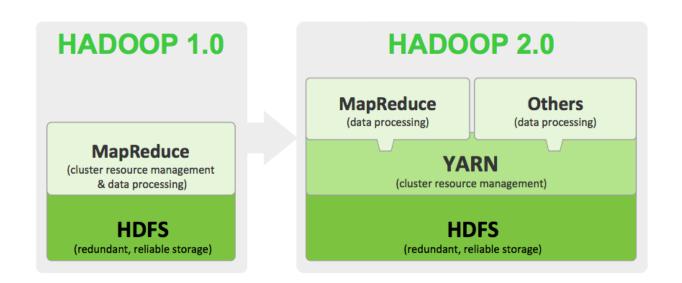




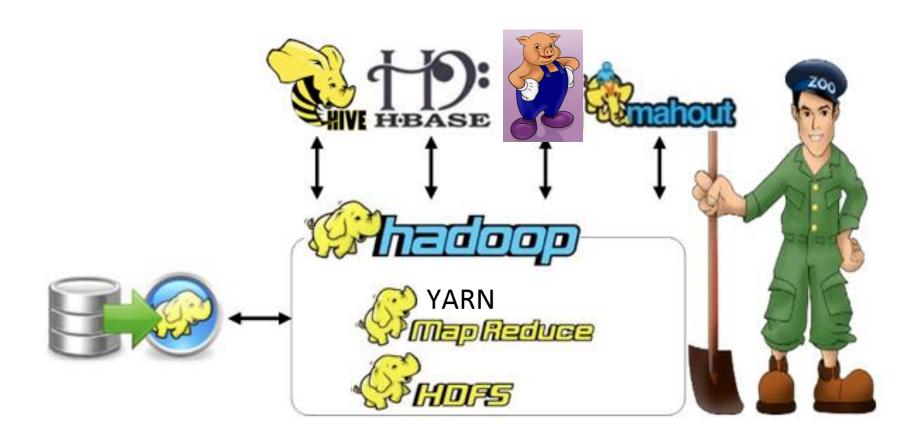
Hadoop 2

YARN

Gestion jobs et ressources



Ecosytème



Plan

- Rappels sur Hadoop
- Présentation (rapide) de Solr
- Intégration d'Hadoop et de Solr
- Démo

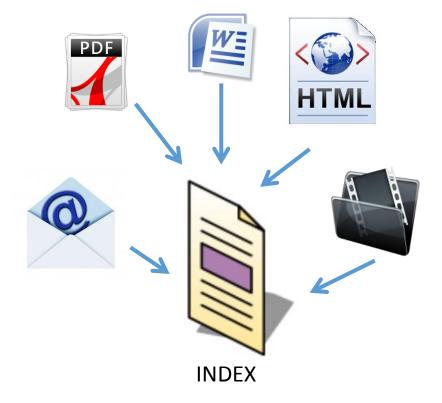


La recherche est un oignon matriciel!

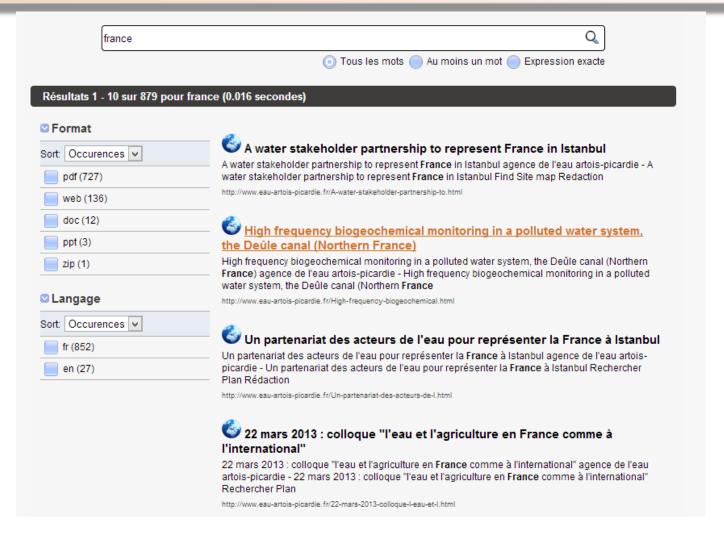


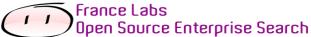
• Un outil qui permet:

De créer un index à partir de documents

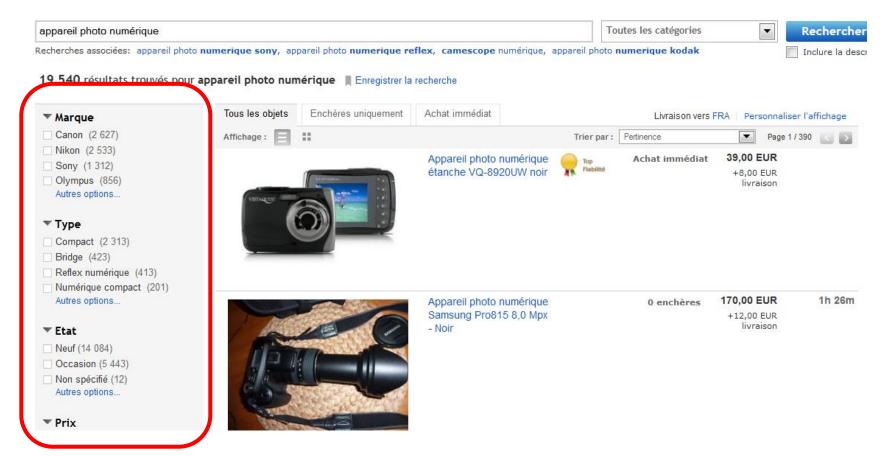


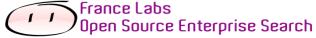
A quoi ça ressemble ?





Facettes





Spellcheck

androit

Related searches for androit: android phone cell phones android tablet android table

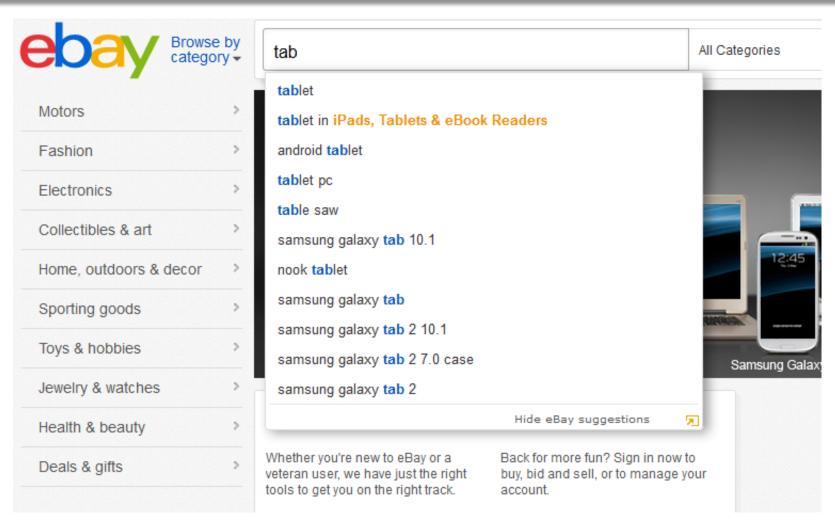
5 active listings | sold listings | completed listings

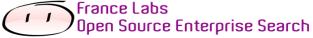


Did you mean: android? (57652 items)

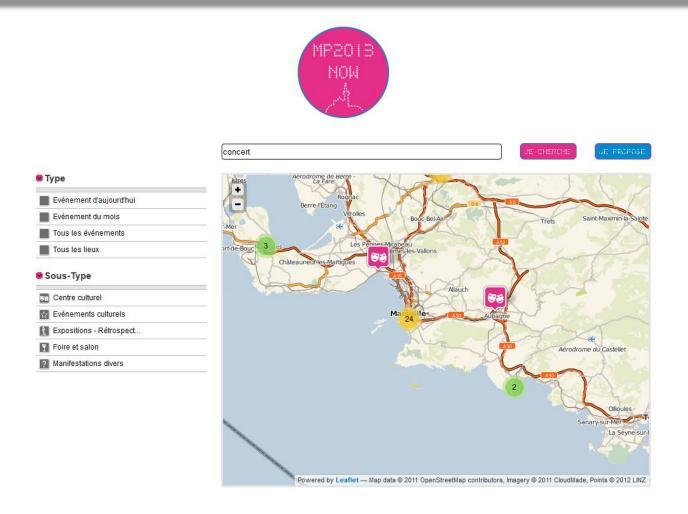


Autocomplete





Géolocalisation





Autres fonctionnalités

Highlighting

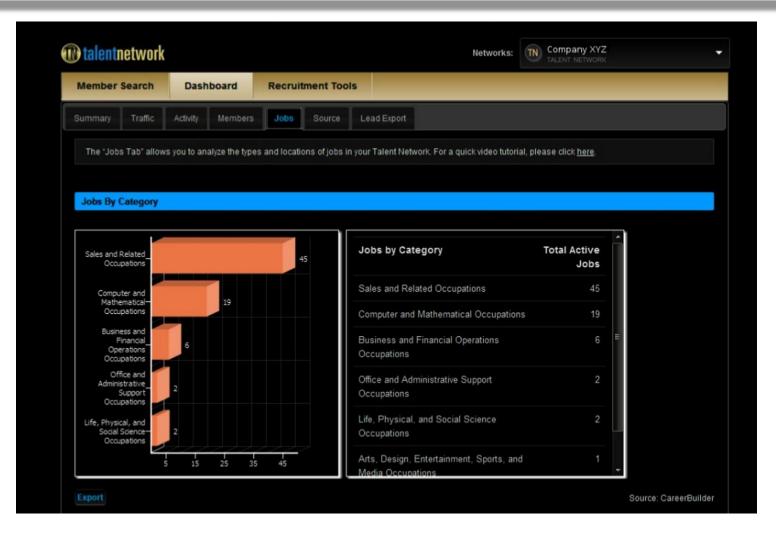
15 janv. 2012 – **Solr** provides a collection of **highlighting** utilities which can be reused by various Request Handlers to include "**highlighted**" matches and ...

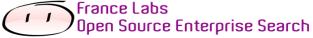
More Like This

- Obtenir des documents similaires à un document
- Similarité textuelle



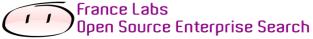
Business Intelligence (BI)/Marketing





Analyse de logs (for devops)



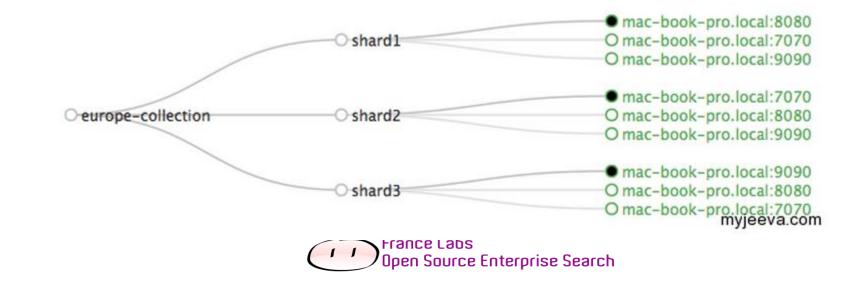


Solr Cloud

Solr Cloud permet la mise à l'échelle:

- Architecture flexible en cluster de serveurs
- Volumétrie élevée
- Montée en charge
- Tolérance aux pannes





Démo Solr

Tester Solr



Démo!



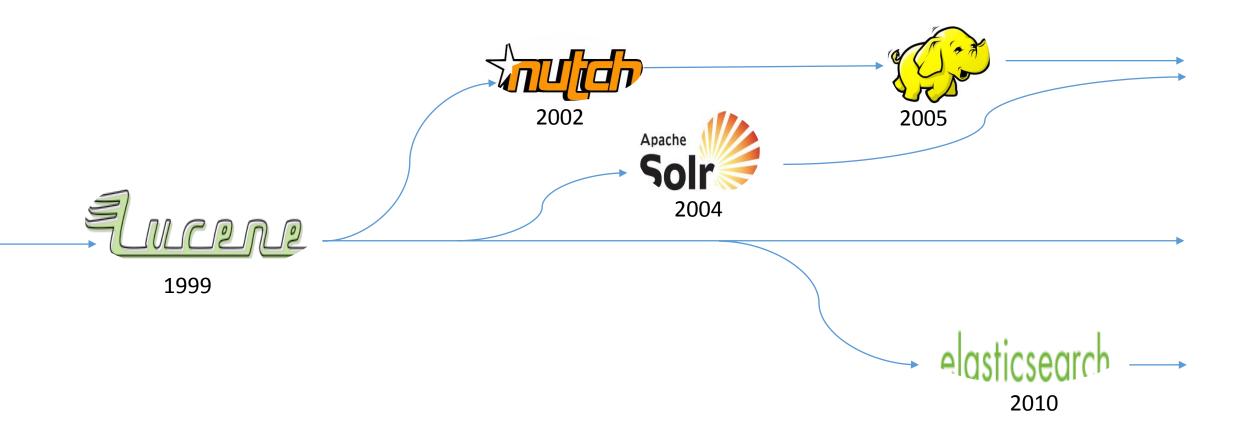


Plan

- Rappels sur Hadoop
- Présentation (rapide) de Solr
- Intégration d'Hadoop et de Solr
- Démo



Projets liés





Exemples de scénarios Big Data

- Text mining
- Création et analyse de graphes
- Reconnaissance de patterns
- Moteurs de recommandations, ciblage publicité
- Analyse de menaces, surveillances de marchés



Hadoop pour les recommandations

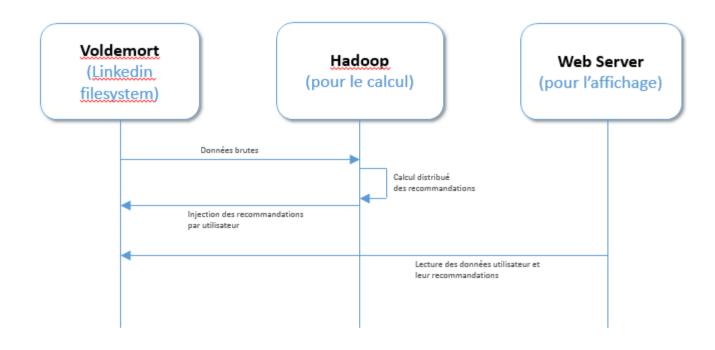


People You May Know





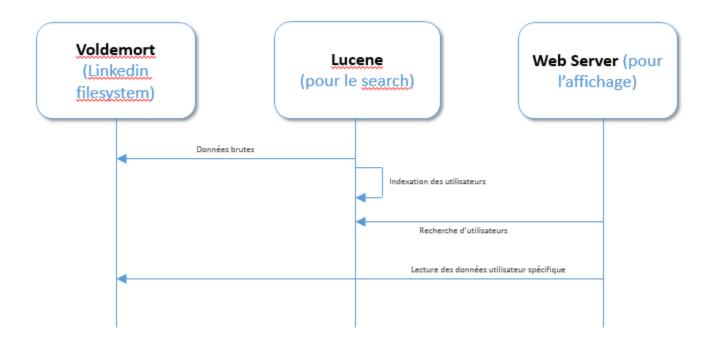
Hadoop pour les recommandations







Lucene pour la recherche d'utilisateurs





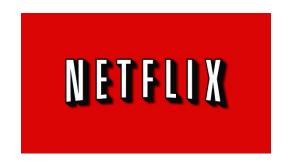


Exemples d'utilisation d'Hadoop

Netflix

- Site de streaming légal de vidéos
- Avant Hadoop :
 - Logs étaient traités pendant la nuit
 - Importés dans une BDD
 - Analyse/BI
 - => Il fallait plus de 24h pour traiter une journée de logs
- Avec Hadoop:
 - En seulement une heure de jobs Hadoop les données sont traitées
 - Traite quotidiennement 1 To de données par jour approximativement
- Solr
 - Recherche de contenus par l'utilisateur





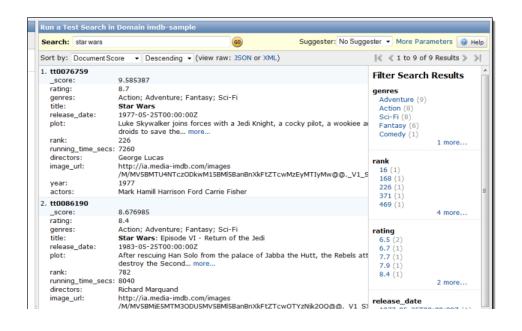


Exemples d'utilisation Big Data

Amazon CloudSearch

- Utilisation de Solr depuis Mars 2014
- Scalable
- Performant

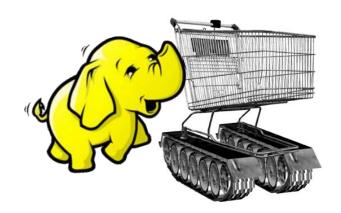






Scénario : analyse de tickets de caisse

- Que faire des données brutes ?
- Comment les analyser, les traiter ?
- Comment rechercher dedans ?
- => Solution : Hadoop et Solr !







1^{er} niveau d'intégration

Hadoop et Solr indépendants

- Traitements et calculs sur Hadoop (ex : Pig)
 - Import d'un CSV avec Pig
 - Opérations : analyse, calculs panier moyen, TVA etc...
- Export des données d'Hadoop
 - Génération d'un fichier XML avec Pig
- Indexation dans un Solr standalone (ex : DIH)
 - Utilisation d'une contribution de Solr : DIH

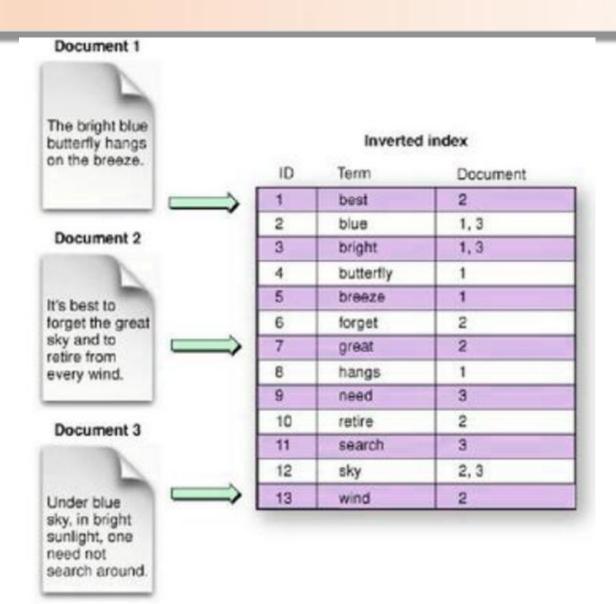








Index Solr



1^{er} niveau d'intégration

Avantages:

- Si cluster Solr déjà présent, ajout d'un index Solr dédié
- Mature

Inconvénients:

 Réplication de HDFS vers système de fichiers du Solr



2^e niveau niveau d'intégration

- Stockage de l'index de Solr sur HDFS
 - Lancement de Solr avec utilisation de HDFS

```
java -Dsolr.directoryFactory=HdfsDirectoryFactory
-Dsolr.lock.type=hdfs
-Dsolr.hdfs.home=hdfs://host:port/path
```







2^e niveau niveau d'intégration

Avantages :

- Tolérance aux pannes
- Gain espace disque (cluster Hadoop déjà existant avec HDFS)









Différents niveau d'intégration

- Création de l'index Solr à l'aide de jobs Map Reduce
 - Travail en cours (https://issues.apache.org/jira/browse/SOLR-1045)
 - Contrib expérimentale fournie avec Solr
 - Avantages :
 - Augmentation des performance d'indexation sur de gros volumes de données
 - Très utile si indexes déjà construits par HDFS





Solr

- Fonctionnalité de Search basée sur Solr intégrée par plusieurs distributions
 - MapR
 - Cloudera Search
 - Hortonworks
 - Lucidworks











Plan

- Rappels sur Hadoop
- Présentation (rapide) de Solr
- Intégration d'Hadoop et de Solr
- Démo



Démo Hadoop et Solr

Work in progress

- Scénario des tickets de caisse
- Etape 1 : import dans Pig des tickets en CSV
- Etape 2 : opérations sur les données avec Pig
- Etape 3 : indexation des données en Solr
- Etape 4 : recherche sur ces données avec Solr







Démo Hadoop et Solr

Work in progress

```
<field name="id" type="uuid" indexed="true" stored="true" required="true"/>
<field name="gtin cd" type="string" indexed="true" stored="true" multiValued="false" />
<field name="gtin nm" type="string" indexed="true" stored="true" multiValued="false" />
<field name="gtin img" type="string" indexed="true" stored="true" multiValued="false" />
<field name="brand nm" type="string" indexed="true" stored="true" multiValued="false" />
<field name="price" type="float" indexed="true" stored="true" multiValued="false" />
<field name="id-show" type="string" indexed="true" stored="true" multiValued="false" />
<field name="shop-name" type="string" indexed="true" stored="true" multiValued="false" />
<field name="categorie" type="string" indexed="true" stored="true" multiValued="false" />
<field name="position" type="location" indexed="true" stored="true"/>
                                                                                          "docs": [
                                                                                              "gtin cd": "0071249812044",
                                                                                              "gtin nm": "Hydra Perfecte Concealer Medium Deep",
                                                                                              "gtin img": "http://product.okfn.org.s3.amazonaws.com/images/gtin/gtin-007/0071249812044.jpg",
                                                                                              "brand nm": "L'Oréal Paris",
                                                                                              "price": 6.21,
                                                                                              "id-show": "shop-48888",
                                                                                              "shop-name": "Champs sur Marne",
                                                                                              "position": "48.85436,2.58183",
                                                                                              "categorie": "article",
                                                                                              "id": "f3bf1047-8b0c-42f7-b77a-1c7d3d48fcd1",
```

" version ": 1466082742653419500

Démo Hadoop et Solr

Solr directement couplé avec Hadoop



Démo!







Contact

Site web: www.francelabs.com

Email: contact@francelabs.com

Twitter: Francelabs

