

# PRESENTATION DU BIG DATA chez ITIM, GTS et BSC

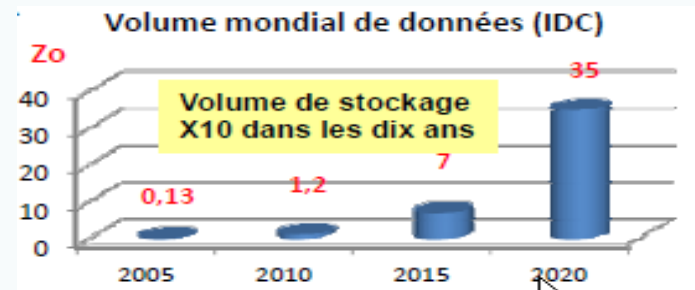


08/09/2017

L'humanité produit aujourd'hui en deux jours autant de données qu'elle en a produites jusqu'en 2003, de toute son histoire.

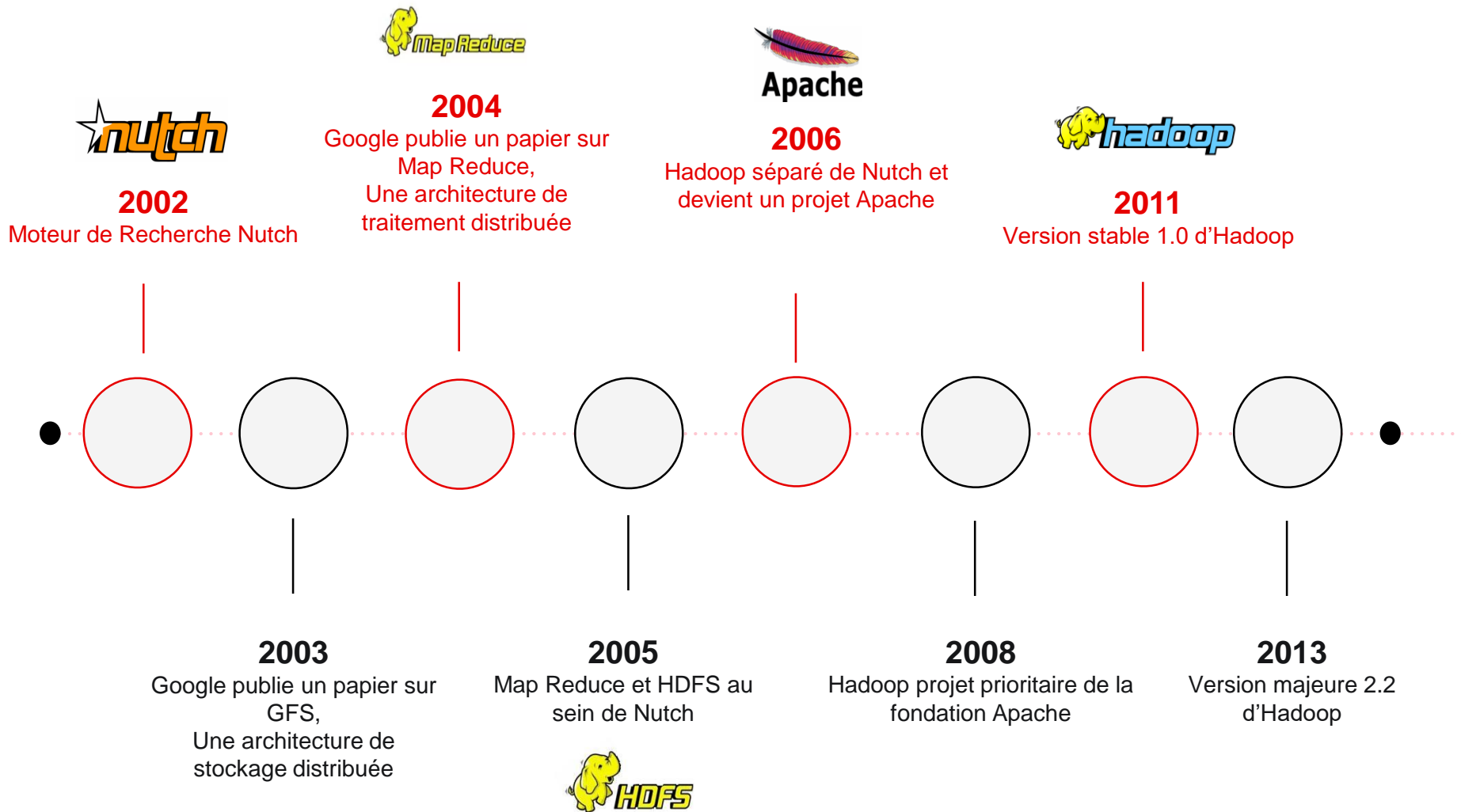


*Explosion des données, structurées, non structurées ou semi-structurées (réseaux sociaux) : textes, dessins vectoriels, images, vidéo, voix, ...*



Gigaoctets ( $10^9$  octets) >> Téraoctets ( $10^{12}$  octets) >> Pétaoctets ( $10^{15}$  octets) >> Exaoctets ( $10^{18}$  octets) >> Zettaoctets ( $10^{21}$  octets)

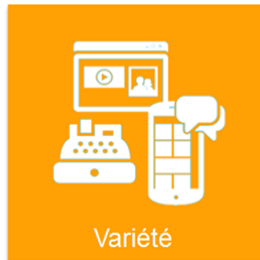
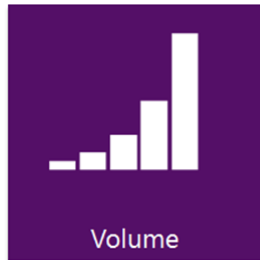
## Un peu d'histoire...



## ➔ Définition du Big Data :

« données structurées ou non dont le très grand volume requiert des outils d'analyse adaptés » (en français : mégadonnées)

## ➔ Caractéristiques du Big Data :



## ➔ Les principaux cas d'usage à ce jour :



### Améliorer la vision 360 du client

En complétant d'informations internes et externes sur, par exemple, le comportement, les centres d'intérêt



### Exploration

Chercher, visualiser et comprendre les données pour améliorer la connaissance du business



### Analyse de données IT

Analyser les données techniques (logs) des applications afin d'améliorer la qualité de service



### Sécurité

Diminuer les risques avec la détection de fraudes et la mise en place d'une « cyber-sécurité » en temps réel.



### Capacité de stockage

Augmenter l'efficacité opérationnelle de l'infrastructure (baisse des coûts de stockage, amélioration des performances, etc.)

	AUJOURD'HUI	DEMAIN
INTERNES	STRUCTUREES  Transactions et opérations bancaires, données de référence, données de synthèse	NON STRUCTUREES  Textes (outils CRM, documents dématérialisés, logs, mails, images, vidéo, ... échangés entre collaborateurs, avec les clients, etc.
EXTERNES	Bases de données externes (partenaires, ...), référentiels publics	Textes, images, vidéos, ... issus d'internet et des réseaux sociaux

En regard du contexte, la construction du Big Data au sein de la Société Générale a adopté les principes suivants :

## DATALAKE



Constituer un « data lake » permettant un stockage unique d'un **maximum de données** (maximum de sources, maximum de variétés) pour répondre à un **maximum de besoins**



## COOPERATION

Faire coopérer les différentes directions de l'entreprise, **facteur clé** de réussite dans la valorisation de la donnée (partage de la donnée)

## SECURISATION



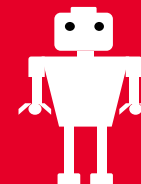
Avoir une politique de **sécurisation interne et externe de la donnée** et des usages pour être conforme aux réglementations sur les données personnelles, et pour **prévenir les risques opérationnels** et notamment le risque d'image

## VISION



Avoir une vision ambitieuse (interconnexion de toutes les données de l'entreprise, voire avec les données externes) et associée à une stratégie des « petits pas » (expérimentation, méthode et structure agile, ...).

## OPEN SOURCE



Etre **compatible** avec les outils Open Source pour être en capacité de suivre les évolutions technologiques

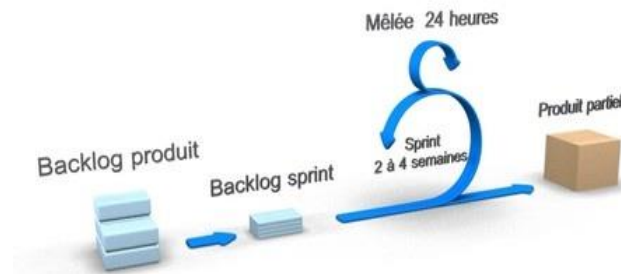
Le contexte nous impose de coller au plus près des besoins des clients, notamment pour réduire le Time to market.



Il est donc nécessaire de proposer un ensemble Big Data « agile » à tous les niveaux :

Agilité dans les **processus** (organisation et méthode)

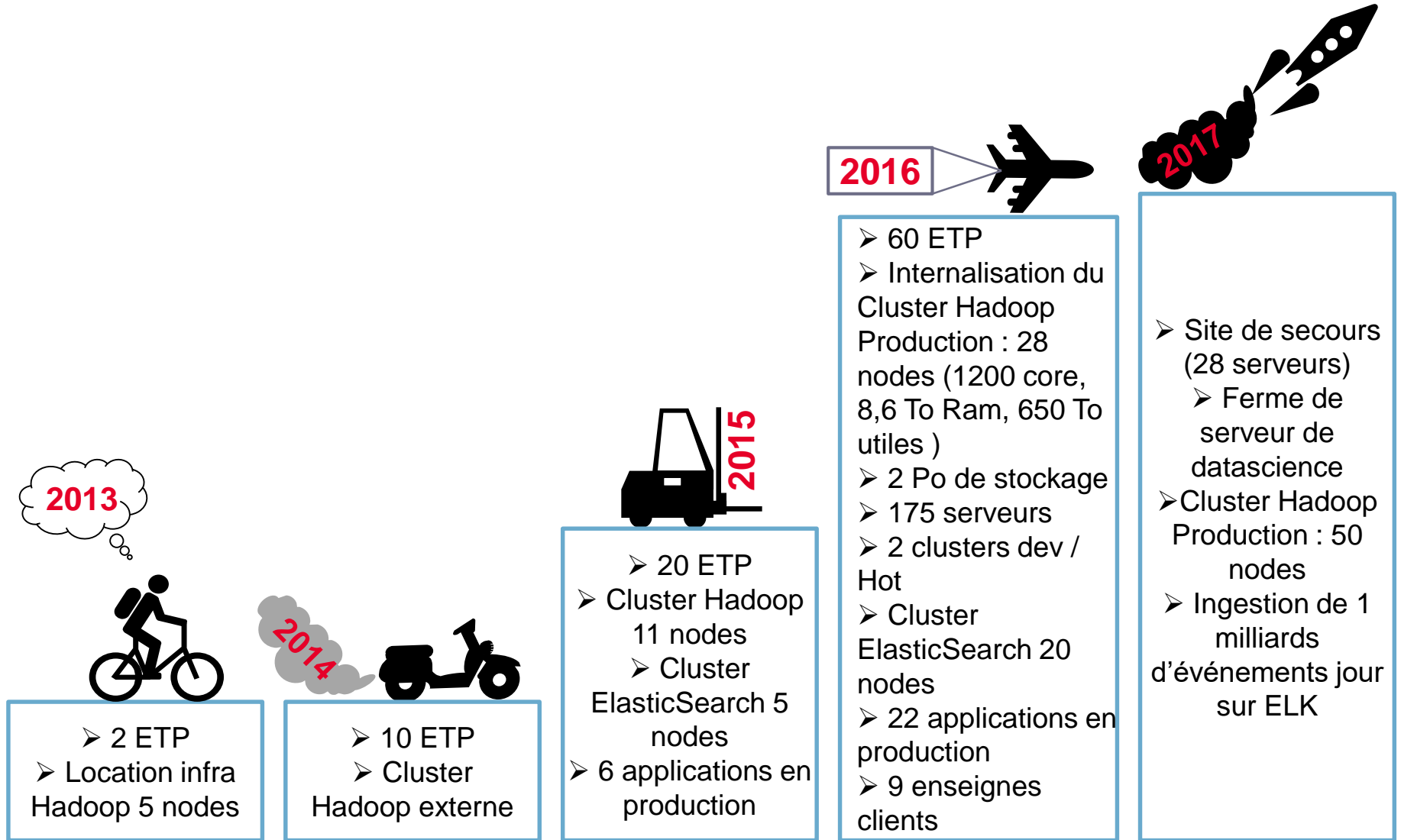
- Organisation : pizza team , proximité avec les métiers, découplage des services
- Méthodes : Lean Startup, développement agile, DevOps

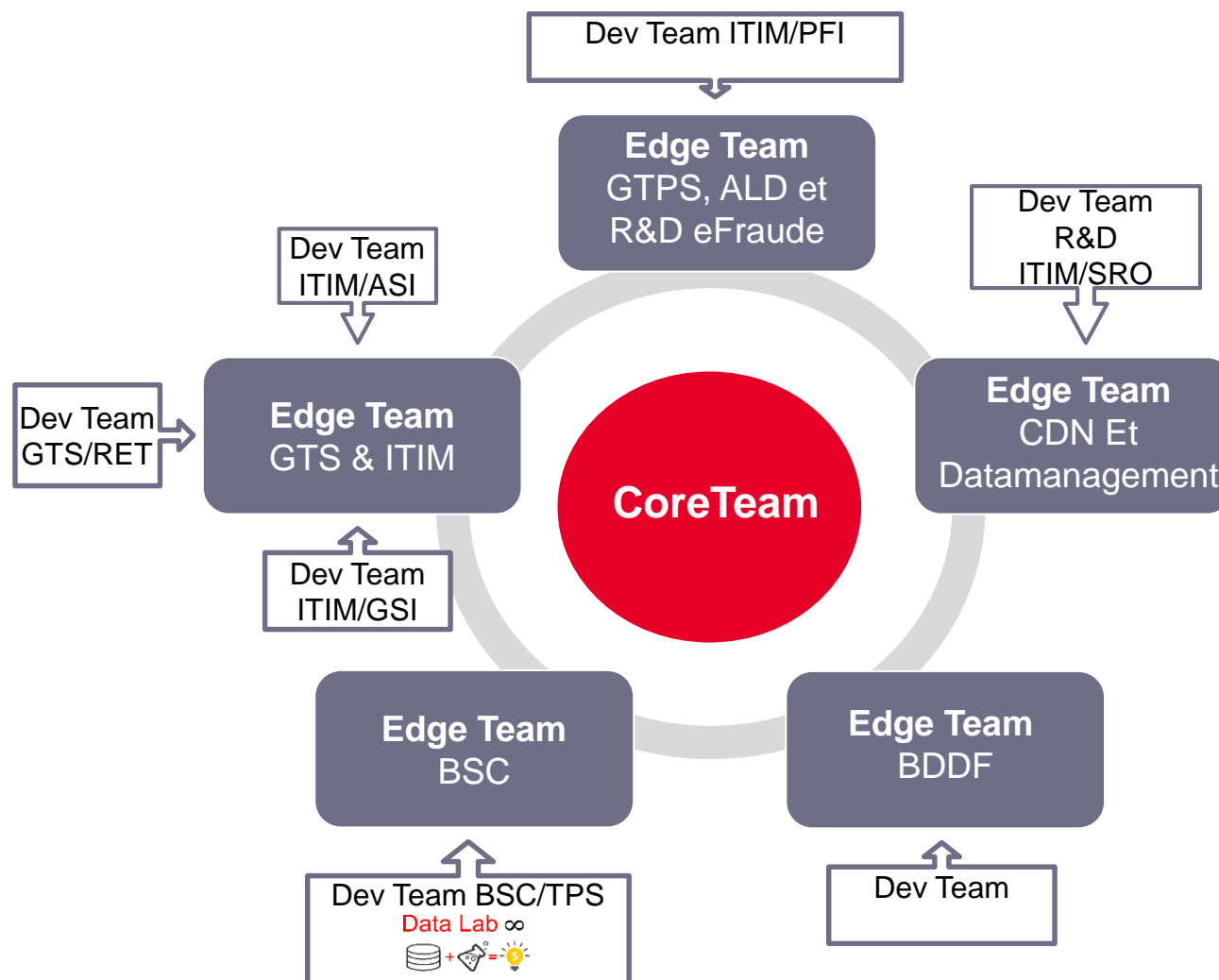


Agilité dans l'**architecture** (logicielle et infrastructure)

- Infrastructure : Construction par itération
- Logicielle : Proximité avec les outils Open Source

# Timeline du Centre de Compétences





## Core Team

Equipe en charge de l'architecture technique, fonctionnelle, du datamanagement, de la sécurité, de l'infrastructure et de l'ensemble des sujets transverses (méthodologie, normes, support)

## Edge Team

Equipe en charge du projet sur des technologies Big Data

## Dev teams

Equipe à l'extérieur du Competency Center en charge d'effectuer les développements à partir des services mis à disposition par la Edge team



## Le Big Data – Domaines d'application

- Amélioration de la connaissance client
- Détection de fraudes
- Développement / Amélioration de scores



EXPLORATION



VISION 360



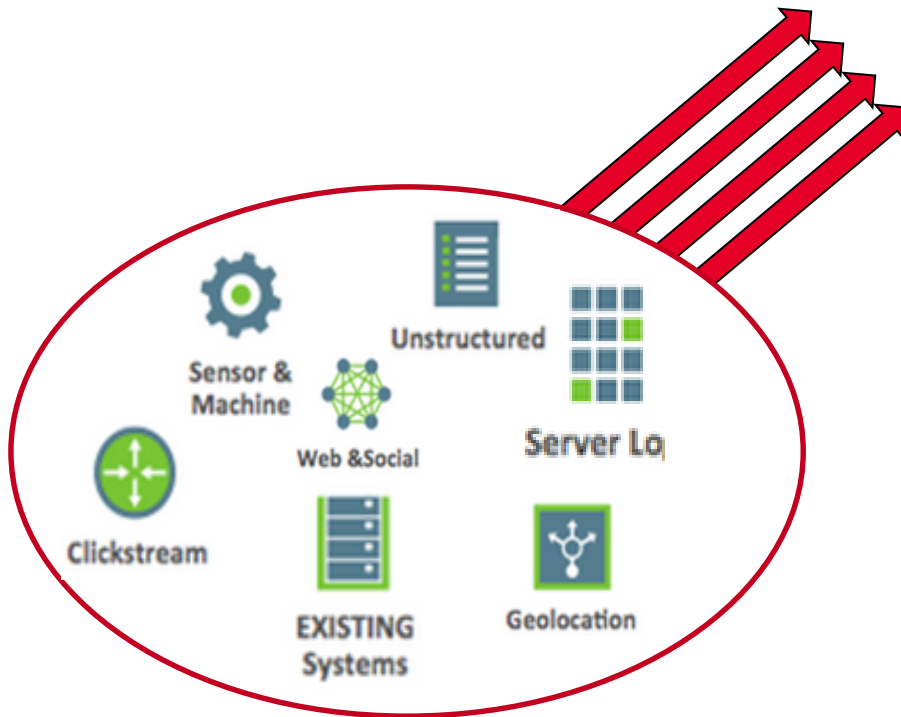
SECURITE



DONNEES IT



STOCKAGE



- Amélioration des process et organisation
- Monitoring et troubleshooting
- Optimisation des temps de traitements

## Supervision & Innovation

« Projet Reporting SHIELD sur les équipements Réseaux Hélios »

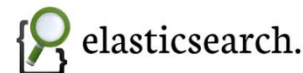
### Enjeux métiers

- Remplacer la plateforme d'hébergement Internet actuelle I2BD par la mise en place d'une nouvelle infrastructure baptisée HELIOS, qui devient en 2015 le nouveau Point d'Accès Internet du Groupe SG
- Répondre aux nouveaux besoins des utilisateurs et à l'accélération des évolutions des nouvelles technologies

### Enjeux SI

- Rendre possible **l'analyse d'un volume important** de données quelques secondes après sa génération par les machines
- **Superviser** en temps-réel les nombreuses technologies mises en œuvre
- **Sécuriser** l'accès aux données, pour les administrateurs comme pour les applications utilisatrices d'Hélios

## Solutions



**30<sup>ème</sup>**  
plateforme  
d'échanges  
avec Internet

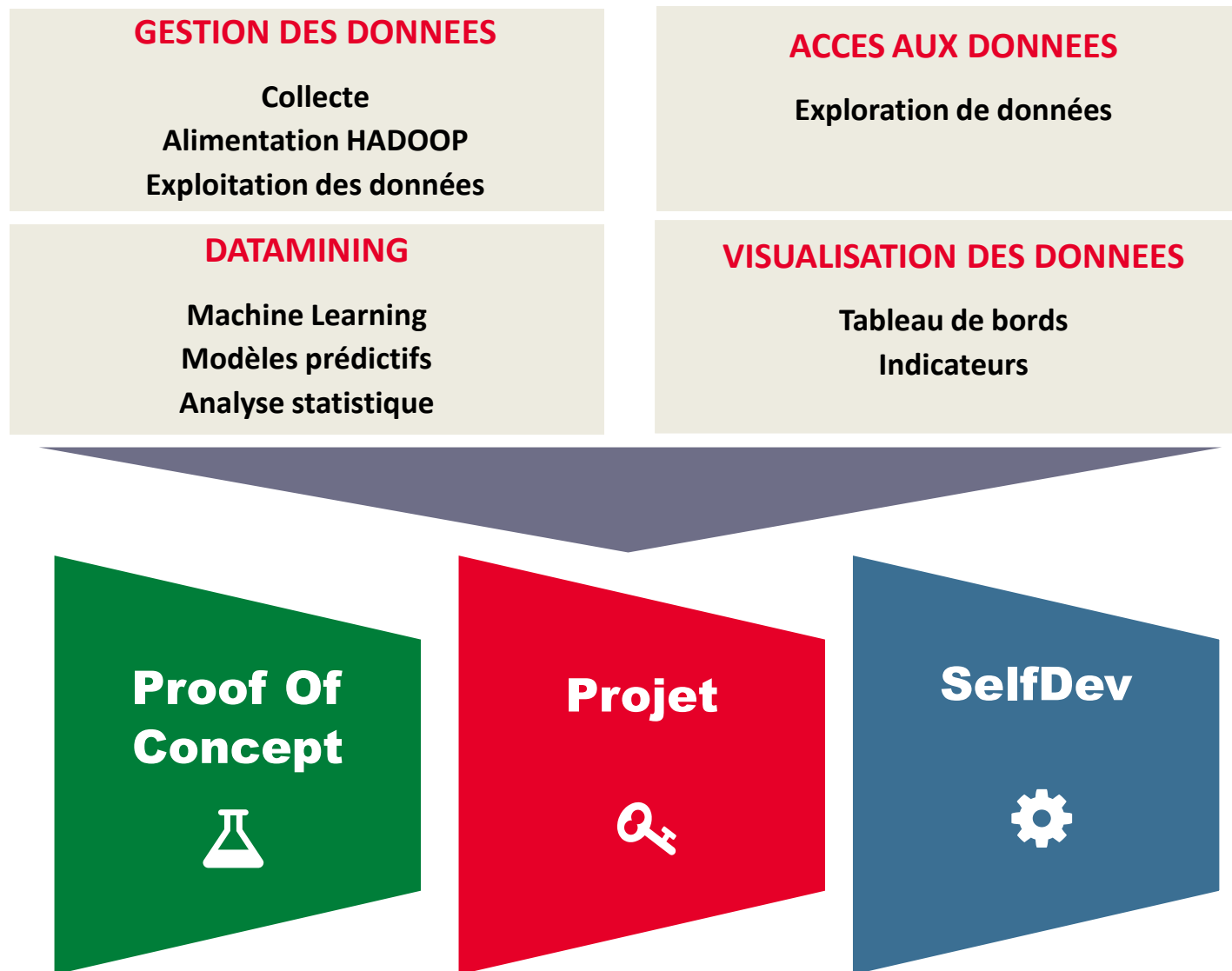
**200** applications  
hébergés en  
France et dans  
le Monde

**9 millions** de  
clients uniques  
par mois

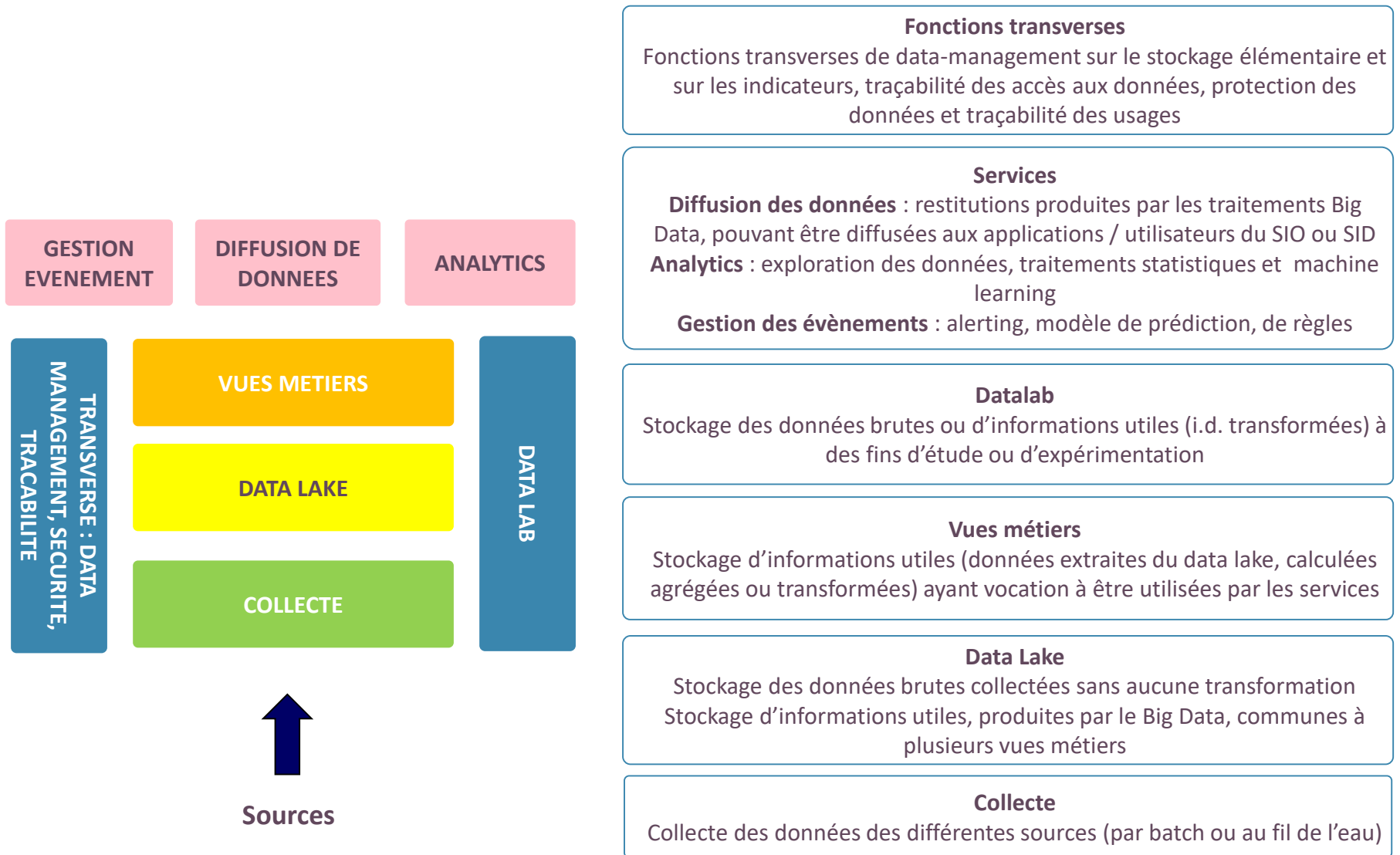
**550** serveurs  
sur 3  
datacenters

**1 milliard**  
d'événements  
par jour  
classifiés C2

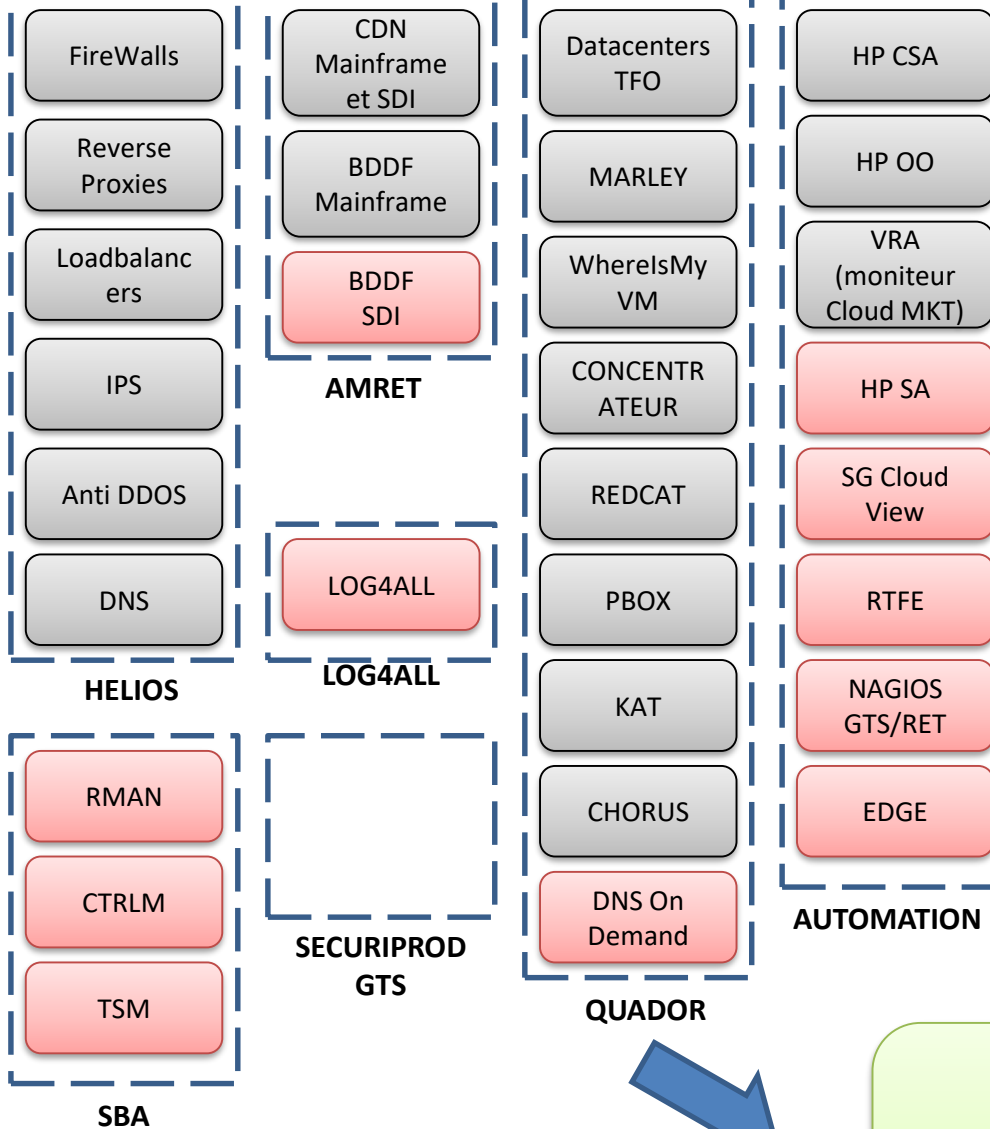
# Offre de services



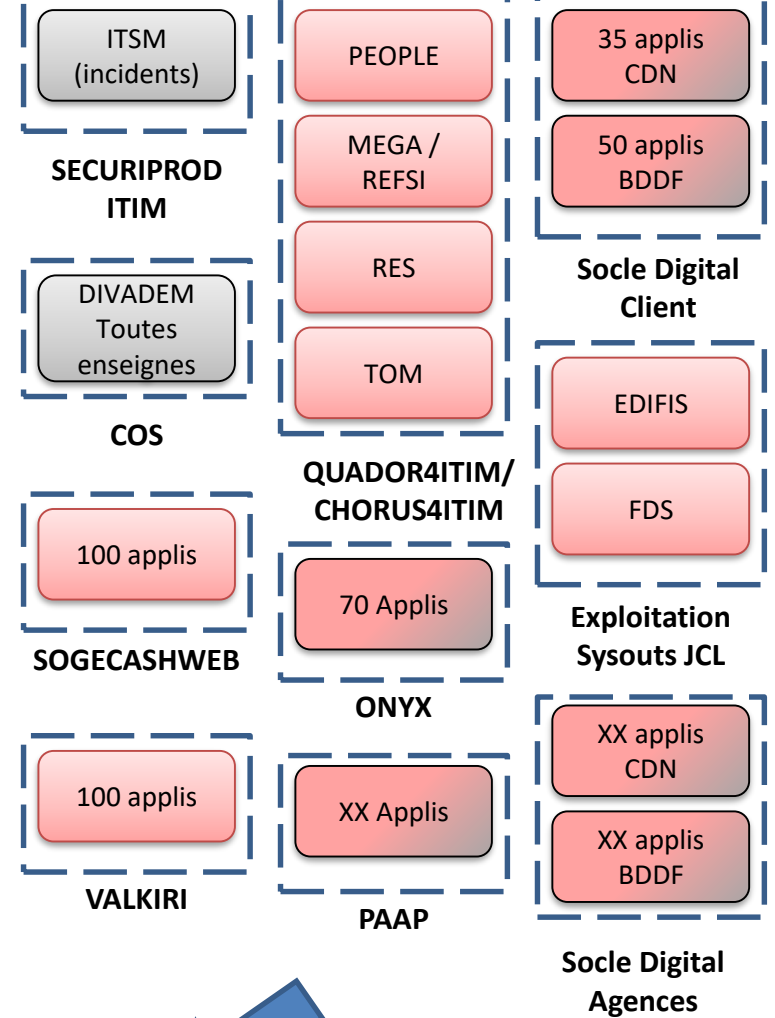
# Architecture des fonctions de la plateforme Big Data



## GTS



## ITIM



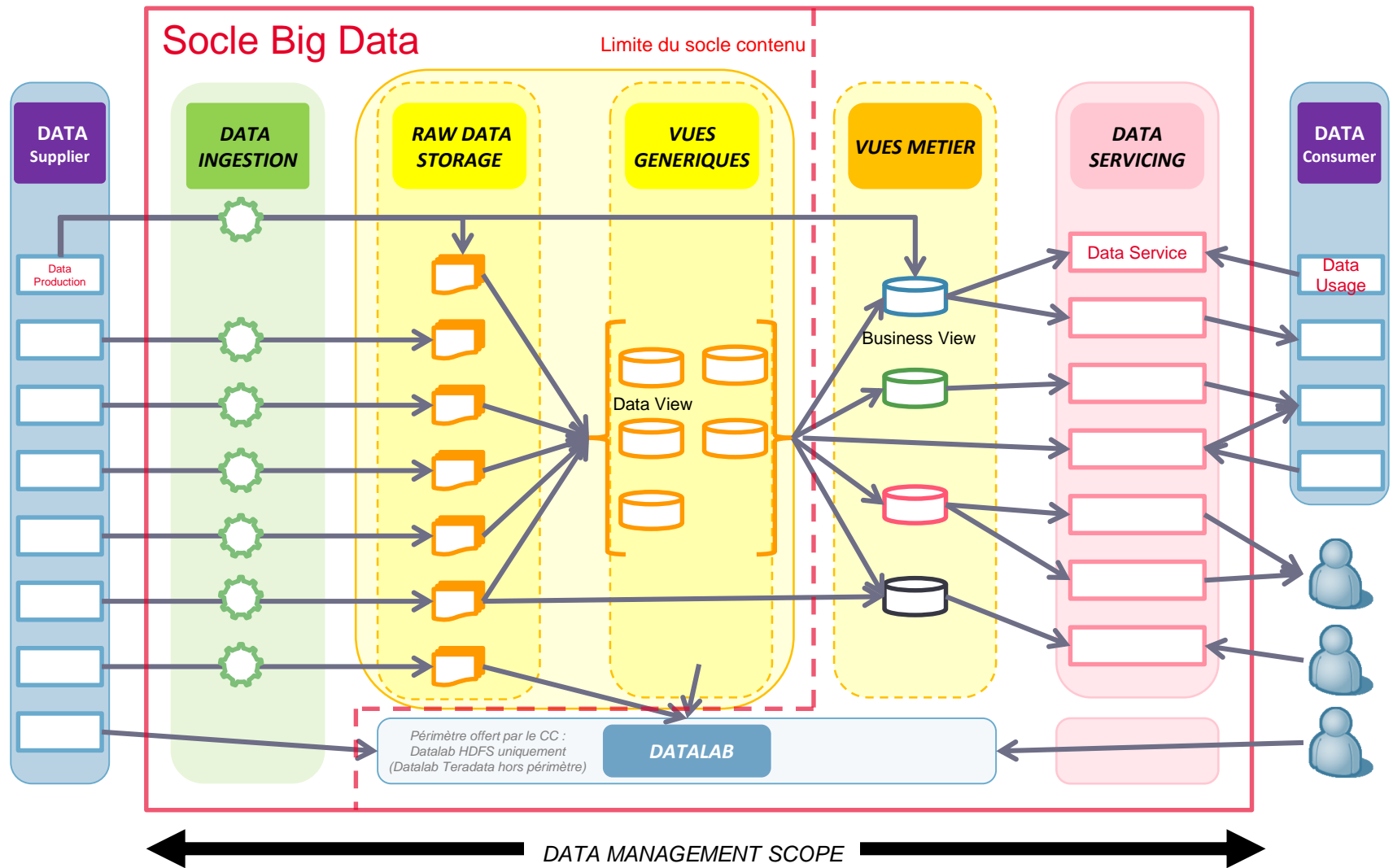
**DATALAKE**

MIND  
GTS

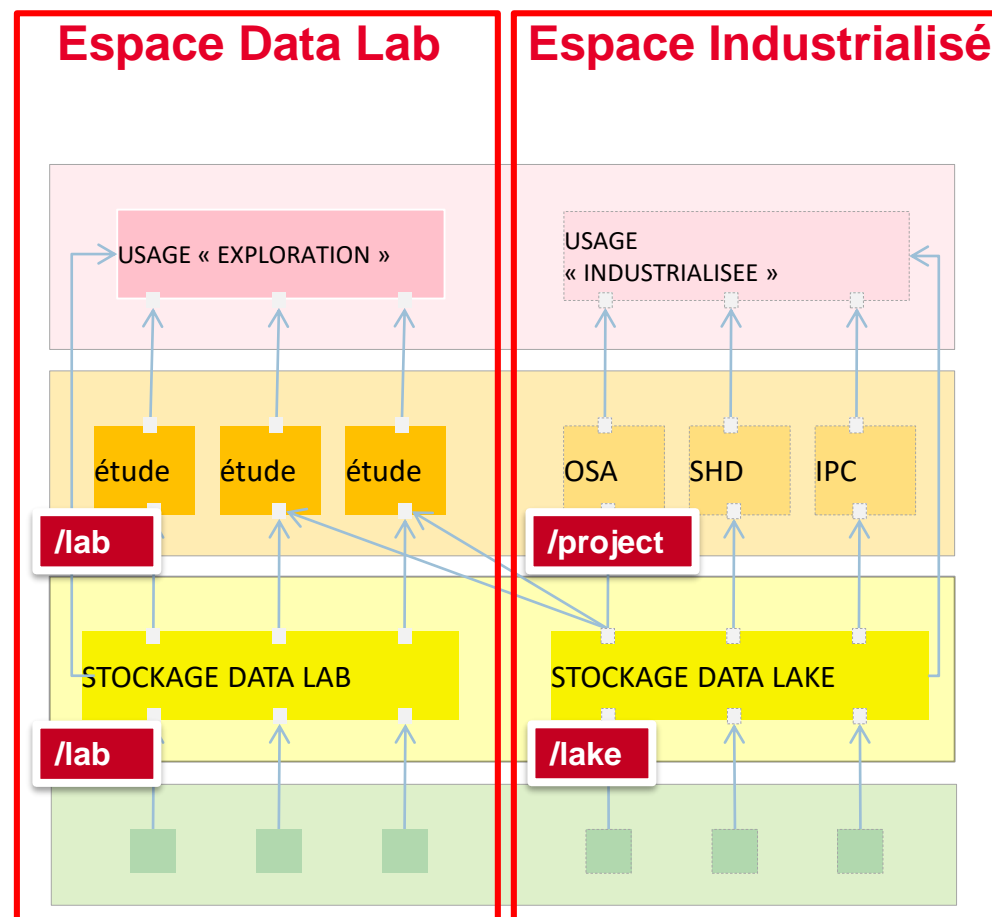
En  
Production

A l'étude

Données  
sortantes



## Exploration des données dans l'environnement Big Data



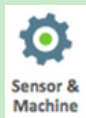
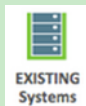
- Usage « exploration » pour le marketing (exemple : manipulation des données avec SAS, R Hadoop, HIVE, Hunk)
- Vues : stockage physique ou logique des données pour un usecase ou un ensemble de usecases cohérents.
  1. Les vues de l'espace Data Lab peuvent accéder aux données de l'espace de stockage Data Lab et Data Lake
  2. Les vues de l'espace industrialisé ne peuvent accéder qu'à l'espace de stockage Data Lake
- Stockage Data Lake : stockage pour les usages métiers industrialisés (process automatisé de la collecte, de calcul d'indicateur, ...) organisé par sources de données
- Stockage Data Lab : Stockage pour les usages métiers adhoc, ie usage oneshot, exploration, recherche de modèle. Les process ne sont pas industrialisés. Le calcul des indicateurs est pris en charge par les métiers (écriture des scripts, lancement, ...). Cf. détail page suivante
- Collecte :
  1. alimentation Data Lab manuelle (via BT ou interface)
  2. alimentation automatisée par batch ou au fil de l'eau

OSA : logs applicatives du poste de travail PUMA  
SHD : logs techniques des points d'accès internet  
IPC : logs de navigation web

# Architecture logicielle

## Des outils sélectionnés et positionnés

### SOURCE DE DONNEES



### INTEGRATION DE DONNEES

#### Ingestion (batch)

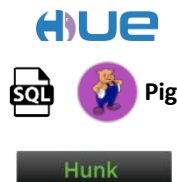


#### Streaming (fil de l'eau)



### USAGES DES DONNEES

#### Analyse batch



#### Analyse interactive



#### Analyse prédictive



#### API



### VUES METIERS

#### Batch



#### Interactif NoSQL



#### Interactif Search Engine



#### Temps réel



### DATA LAKE - DATALAB

HDFS : système de fichiers distribué



### TRANSVERSE : DATA MANAGEMENT, SECURITE, TRACABILITE

Data Management

Habilitations : Ranger, Shield, LDAP, Knox

Audit : Ranger

Cryptage : hp Encryption

### EXPLOITATION

Supervision : Ambari, Nagios, Ganglia, Marvel

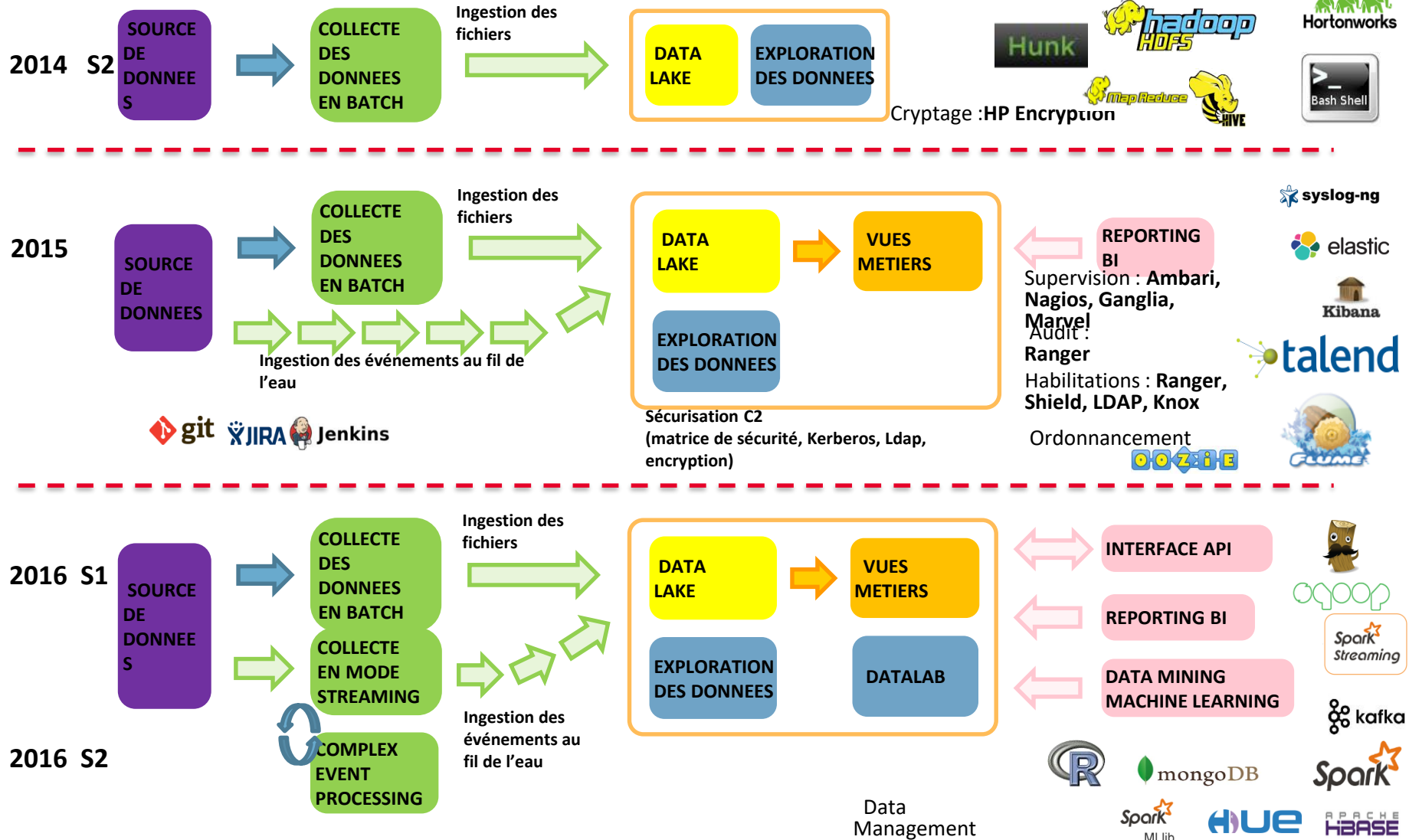
Ordonnancement  
OOZIE

### DEVELOPPEMENT

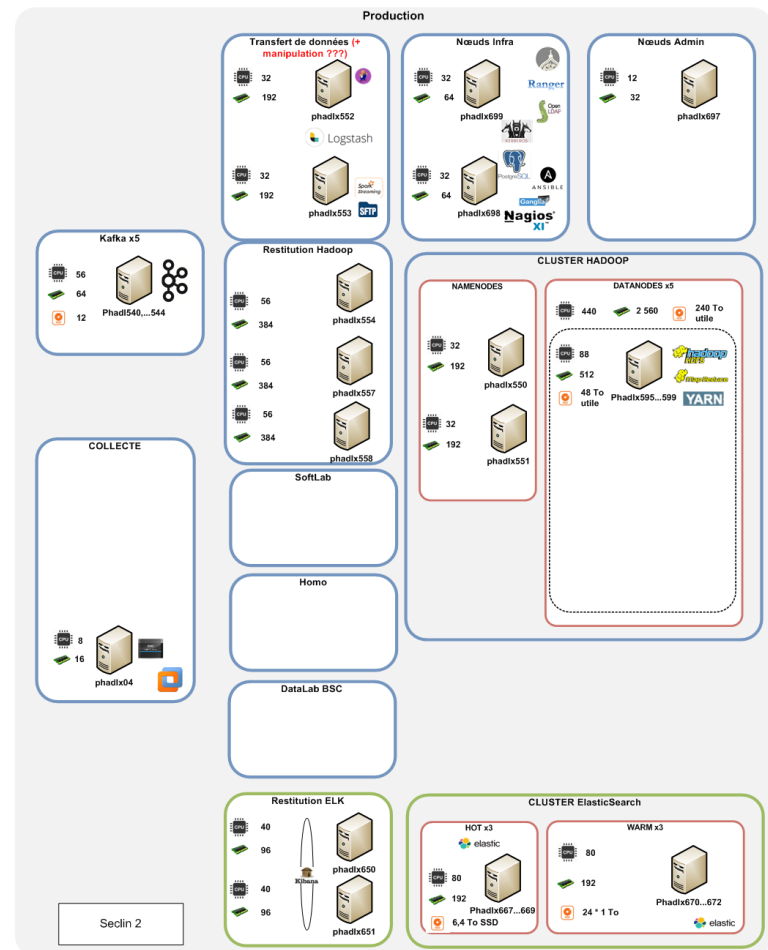
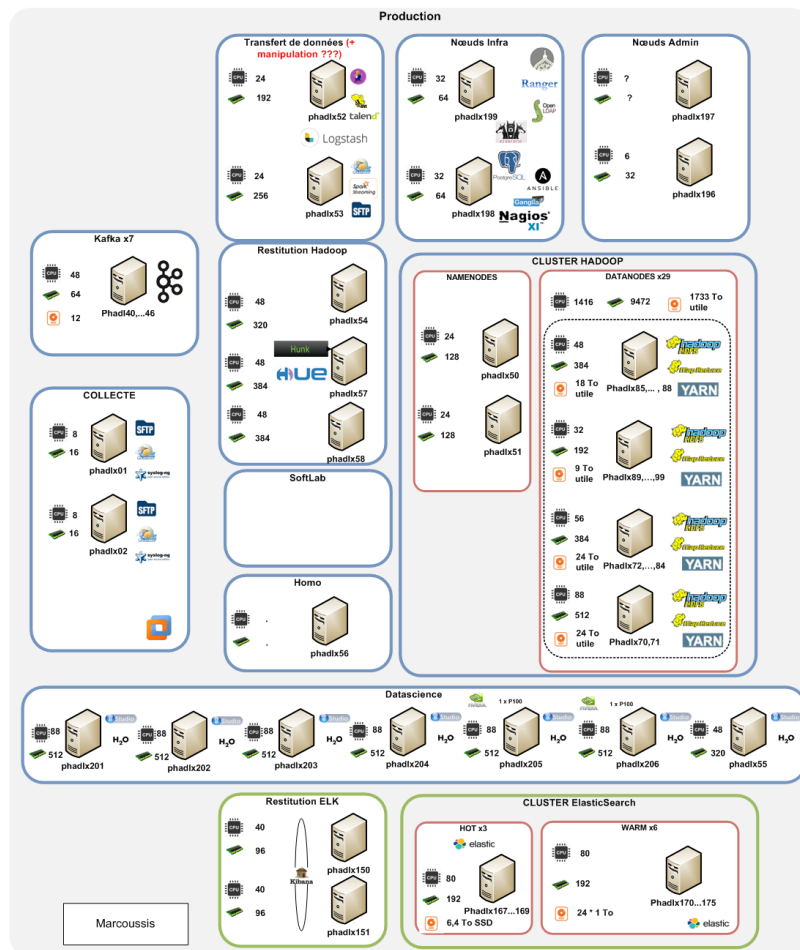
talend, DEV OPS, git, JIRA, Jenkins, Nexus



# Le BIG DATA – L'intégration dans le SI au travers des travaux du CC.



# Architecture physique



- **Mettre en place des outils de protection de la données de bout en bout**

- Chiffrement global du transport des données sur toute l'infrastructure
- Chiffrement des données stockées, sur le NAS GTS pour le serveur de collecte, et sur disque local chez HP avec transmission des clés depuis 2 appliances HSM situées chez GTS

- **Renforcer la traçabilité et les contrôles de tous les composants de la solution**

- Traces d'audit activées sur l'ensemble de l'infrastructure (syslog et auditd) et archivage chez LogColCor. Traçabilité des actions d'exploitation GTS via l'outil ObservIT.
- Supervision des logs via des rapports et alertes HUNK et bientôt Log4All

- **Renforcer la sécurité des serveurs exposés**

- durcissement OS et détection des vulnérabilités par le SOC

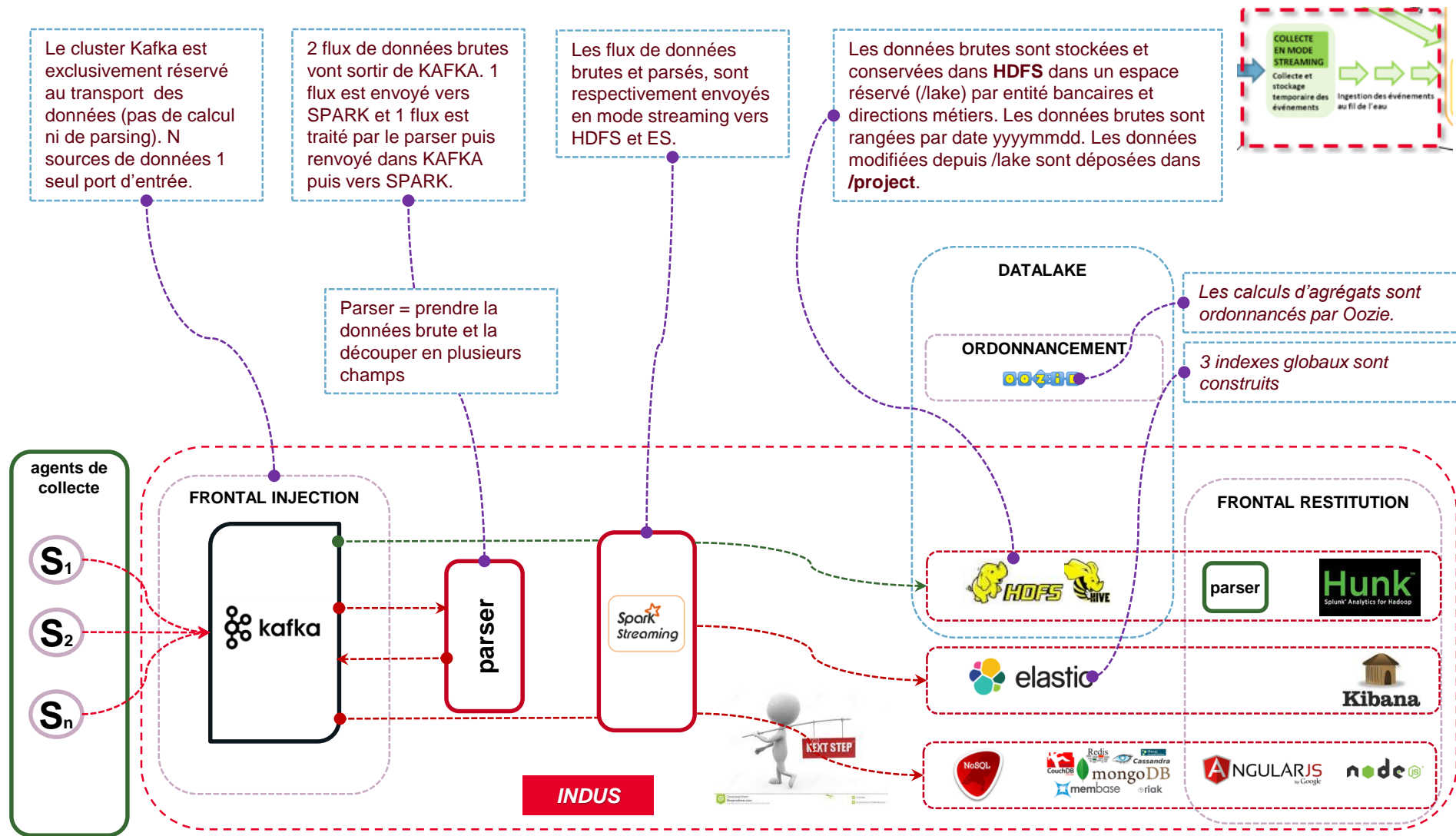
- **Renforcer l'authentification et les habilitations**

- Gestion des comptes utilisateurs depuis SAFE
- Retrait des comptes à privilèges chez l'exploitant (dont root) au profit d'un catalogue de commandes prédéfinies
- Solution automatisée d'export des bases de comptes (LDAP, SAFE) pour revue mensuelle des droits
- Sécurisation des accès au cluster : ACL (Ldap, Ranger, Knox, Shield pour ELK) et Kerberos

- **Mettre en œuvre une démarche risque adaptée pour les nouveaux cas d'usage**

- Audit Sécurité à chaque livraison d'un dispositif Sécurité

# Principe d'alimentation en streaming

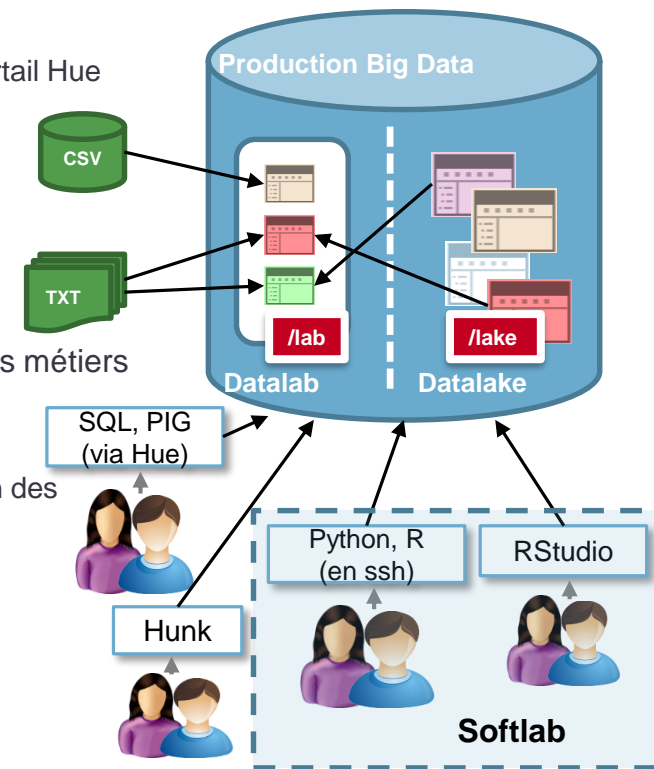


### Datalab :

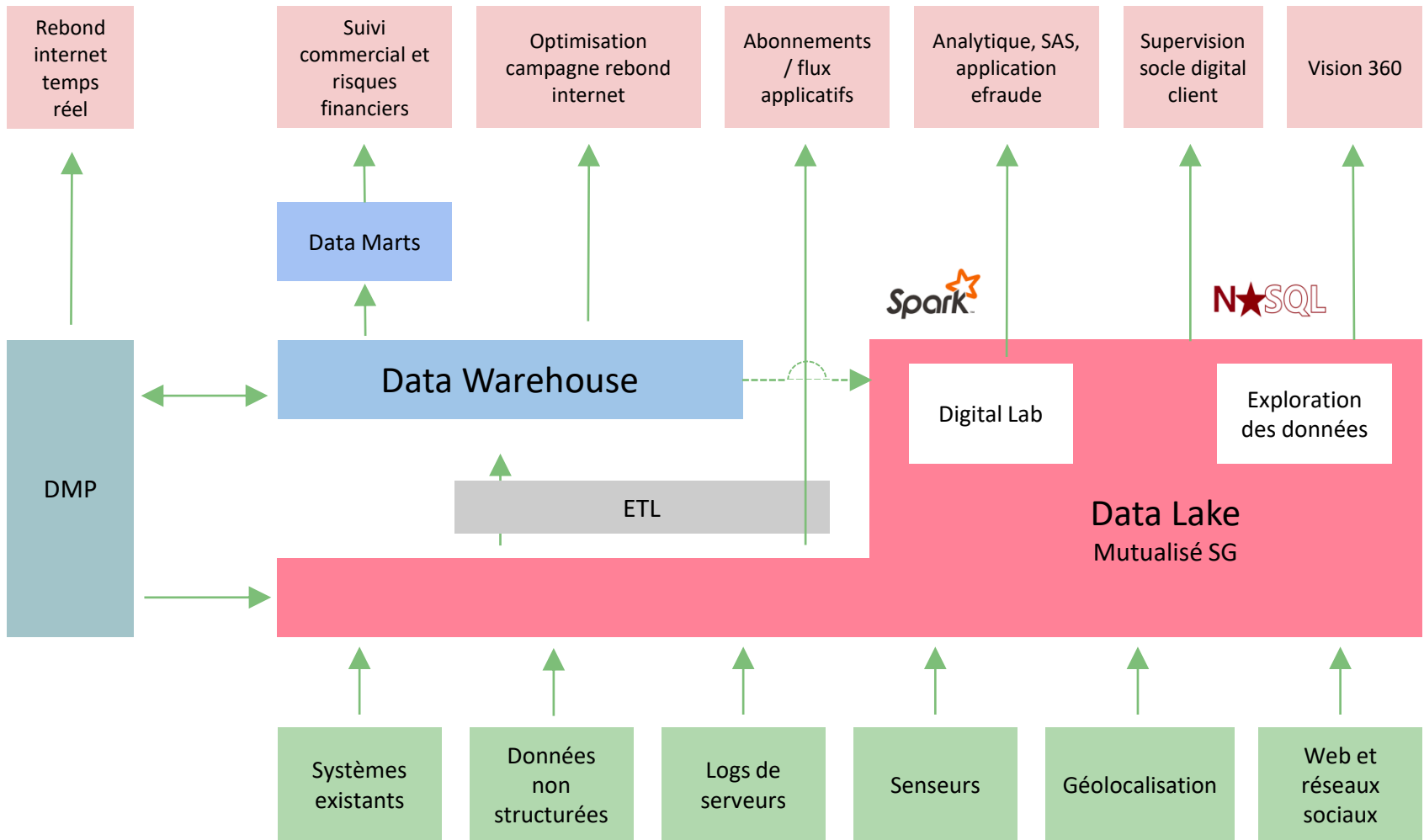
- **Objectif** : offrir un espace de travail en production pour des besoins d'étude, de découverte de données, de modélisation ou de prototypage
- **Principes** :
  - Espaces privatifs par direction métier (lecture / écriture, quota par direction)
  - Chargement et utilisation possible de données externes
  - Accès en lecture aux données du datalake
  - Avec les performances de la production mais non prioritaire par rapport aux traitements industriels et sans automatisation des traitements
  - Accès aux données :
    - en utilisant Hive SQL, langage PIG, navigateur de fichier depuis le portail Hue
    - en utilisant des outils d'analyse (Hunk)
    - en utilisant des outils du Softlab
- **Mise à disposition des datalabs** :
  - création des datalabs effectuée pour les équipes BSC, BDDF, CDN et SRO

### Softlab :








- **Objectif** : proposer un dispositif en production permettant de mettre à disposition des métiers des outils dans le but de les tester ou de les utiliser dans un mode non industriel
- **Principes** :
  - Choix des outils par le métier avec une évaluation du CC Big Data (licence, sécurité, gestion des ressources). Le métier se charge de la relation avec le fournisseur.
  - Installation des outils par le CC Big Data mais pas de support applicatif
  - Traitements lancés par le Softlab moins prioritaires que les traitements industriels
  - Réinstallation possible des outils du Softlab sur un serveur industriel au besoin
- **Mise à disposition du Softlab** :
  - mise à disposition en production du serveur avec installation des outils :
    - Python 3.4 et ses différentes librairies
    - R 3.2.2, ses différents librairies ainsi que l'interface web Rstudio Server
    - configuration pour permettre les accès en ligne de commande (ssh)



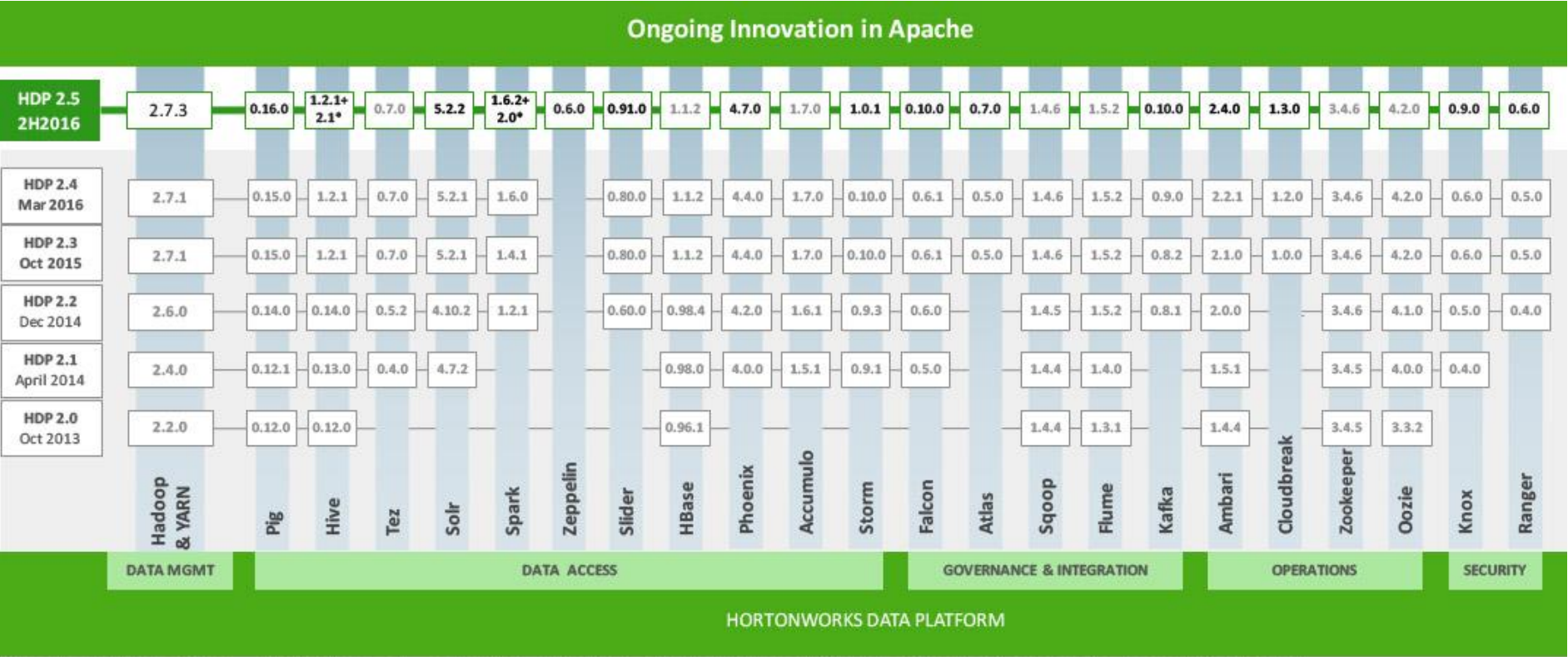
## Architecture de référence détaillée – Schéma général



## Annexe 1 : Choix des outils d'ingestion/diffusion

Thèmes	Sous thèmes	Avantage	Use cases	Produits						
										
Ingestion Traitement Batch	ELT (Extract Load Transform)	<ul style="list-style-type: none"> <li>Scalabilité</li> <li>Traçabilité</li> <li>Gain en maintenance</li> </ul>	Tous types de use cases	✓						
	Spécifique s'exécutant sur l'ensemble du cluster	<ul style="list-style-type: none"> <li>Scalabilité</li> <li>Optimisé pour le use case</li> </ul>	Pour des besoins particuliers nécessitant des optimisations			✓				
	ETL	<ul style="list-style-type: none"> <li>Traçabilité</li> <li>Gain en maintenance</li> </ul>	Pas recommandé		?					
Ingestion Traitement fil de l'eau	Collecte et ingestion	<ul style="list-style-type: none"> <li>Résilience</li> <li>Scalabilité</li> </ul>	Socle digital (IED)			✓	✓		✓	?
	Ingestion	<ul style="list-style-type: none"> <li>Scalabilité</li> </ul>	Hélios (HLS)					✓	✓	
Complex Event Processing	Calcul événementiel	<ul style="list-style-type: none"> <li>Résultat de calcul au plus près du temps réel</li> </ul>	eFraude (EFR)			✓				

## Annexe 2 : Distribution Hortonworks



\* Spark 1.6.2+ Spark 2.0 – HDP 2.5 support installation of both Spark 1.6.2 and Spark 2.0. Spark 2.0 is Technical Preview within HDP 2.5.  
Hive 1.2.1+ Hive 2.1 – Hive 2.1 is Technical Preview within HDP 2.5.



## Annexe 3 : Principe de l'architecture technique

### Serveur d'accès client pour les accès typés industriels :

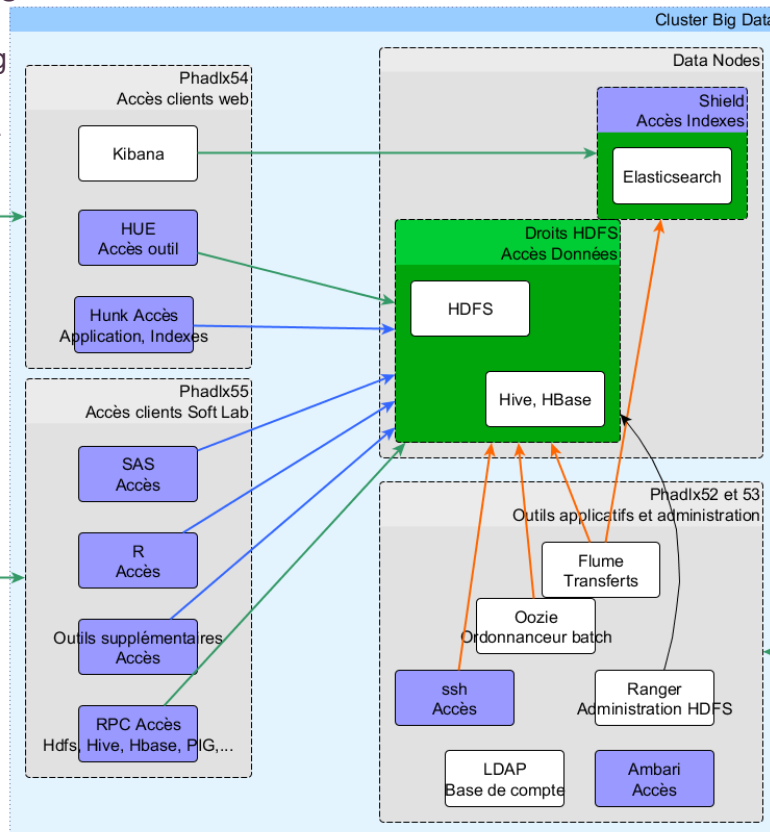
- interface Kibana pour ElasticSearch
- interface Hunk pour exploration, visualisation et reporting Hive, HDFS
- interface Hue pour manipuler les fichiers, lancer des programmes Pig et du SQL Hive
- frontal d'accès pour les outils SQL type MicroStrategy et BO

Usage "Industrialisé"  
Accès Clients industriels

Usage "Data Lab"  
Accès Métier  
et équipes de développement

### Serveur d'accès client pour les besoins d'expérimentations :

- accès SAS
- accès R Hadoop
- outils complémentaires à tester
- accès unix client (non disponible pour l'instant)



### DataNode :

- serveur de traitement et de stockage
- comprenant des espaces Data Lab et industrialisé
- allocation de ressources dynamiques en fonction d'un profil de consommation de ressources (queue YARN)

- Groupe Rôle** - sécurité d'accès aux données géré par des groupes
- Groupe Data** - données organisées par enseigne/direction pour le Data Lab et par SI/application source pour l'espace industriel



### Serveurs d'accès pour les besoins d'administrations et pour l'automatisation des traitements :

- interface d'administration du cluster : Ambari
- outillage de planification et d'automatisation
- serveurs en haute disponibilité pour assurer les alimentations industrielles en cas de problème ainsi que l'administration

## Annexe 4 : CoPil d'architecture Big Data

### Executive Summary

#### • Notre situation actuelle

- Nous disposons d'un Data Lake et de DataWarehouse indépendants
- Nous souhaitons intégrer le Data Lake dans l'architecture du SID (Système d'Information Décisionnel) afin de promouvoir les usages exploratoires
- Nous avons étudié la façon d'intégrer le Data Lake et nous avons réfléchi aux principes de construction et aux impacts de l'architecture de référence

#### • Les principes de notre architecture de référence

- **Agilité** : Privilégier les usages exploratoires grâce à la flexibilité permise par le Data Lake
- **Ouverture de la donnée** : Mise à disposition rapide de la donnée au niveau le plus granulaire possible
- **Utilisation pragmatique** : Choix des solutions les plus adaptées à chaque usage
- **Rationalisation et simplification des SI** : Centralisation progressive de la donnée dans le Data Lake et simplification de l'utilisation des infrastructures pour l'utilisateur final
- **Anticipation** : Préparation aux usages futurs et amélioration du time-to-market
- **Pérennité** : Mise en place de règles d'utilisation afin de maintenir l'ordre sur le Data Lake et de garder la maîtrise des données

#### • Impacts de l'intégration du Data Lake

##### • Conclusions

- Le Data Lake peut servir de solution de stockage à bas coût grâce à l'hybridation des types de stockage permise par Hortonworks
- L'usage du Data Lake est sécurisé, et sa résilience est en cours d'amélioration
- De nombreuses actions visant la mise en place de règles de gouvernance de la donnée sont en cours

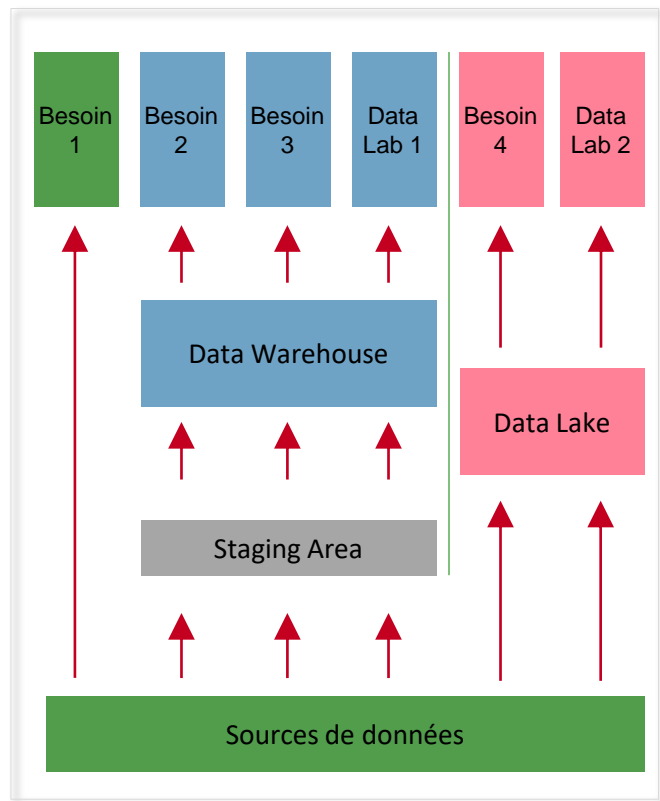
##### • Orientations à valider

- Validons nous le fait de baser notre **système d'alimentation au fil de l'eau** sur les **capacités Big Data** ?
- Devons-nous **restreindre la quantité d'outils** qui seront utilisés en sortie du Data Lake, notamment au niveau des **usages exploratoires** ?
- Partageons-nous l'idée de mettre en place des **grilles d'outils recommandés** pour chaque usage sur la base d'études internes ?
- Devons-nous mettre en place des **outils de gouvernance / Data lineage** de la donnée ?
- Validons-nous le fait de **rapatrier les flux de données existants** vers le Data Lake de **façon progressive**, au fur et à mesure de la **présence avérée de use cases** ?

## Annexe 4 : CoPil d'architecture Big Data

Nous avons alors repensé l'architecture afin de profiter des avantages du Data Lake

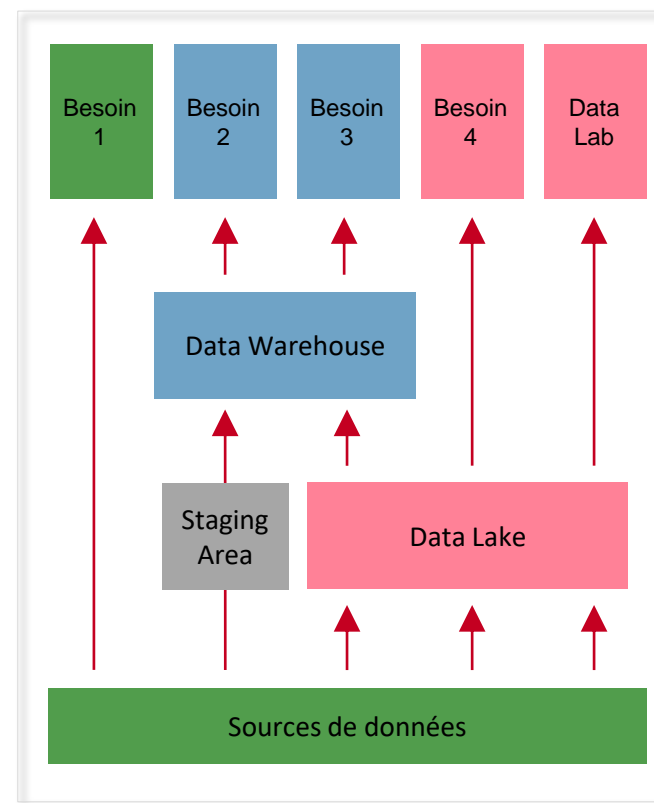
Architecture existante



### Principales évolutions :

- ✓ Levée de la séparation en **silos**
- ✓ **Intégration cohérente** du Data Lake plus qu'une fusion
- ✓ Utilisation du Data Lake pour le **Data Staging**
- ✓ Concentration des **usages exploratoire** sur le Data Lake
- ✓ Construction d'un système d'**alimentation au fil de l'eau**
- ✓ Mise en place de processus de **gestion de la donnée**

Architecture de référence \*



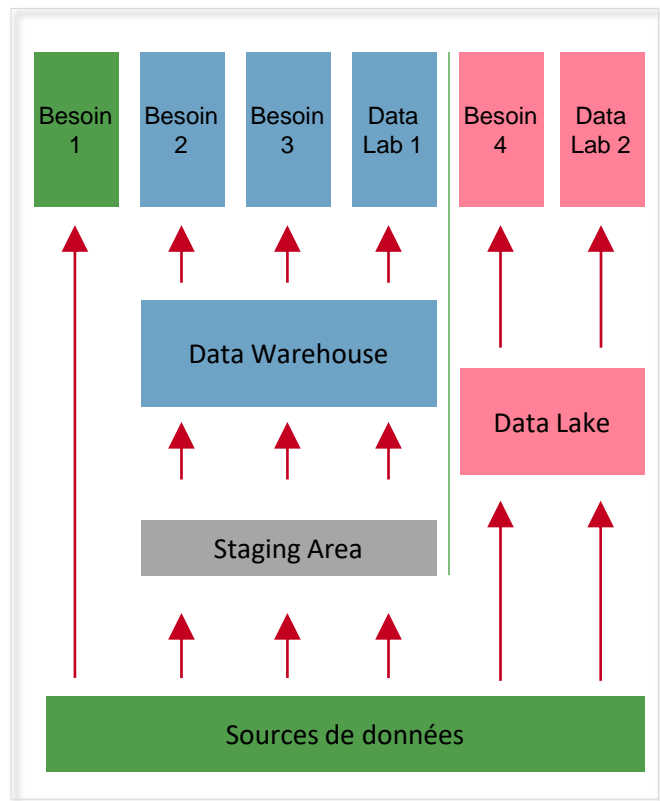
\* Le schéma général d'architecture de référence et son instantiation pour BDDF et CDN se trouvent en annexe.

## Annexe 4 : CoPil d'architecture Big Data

Les évolutions vers la nouvelle architecture seront mises en place de façon progressive

Les Data Labs  
seront basés  
uniquement sur le  
Data Lake

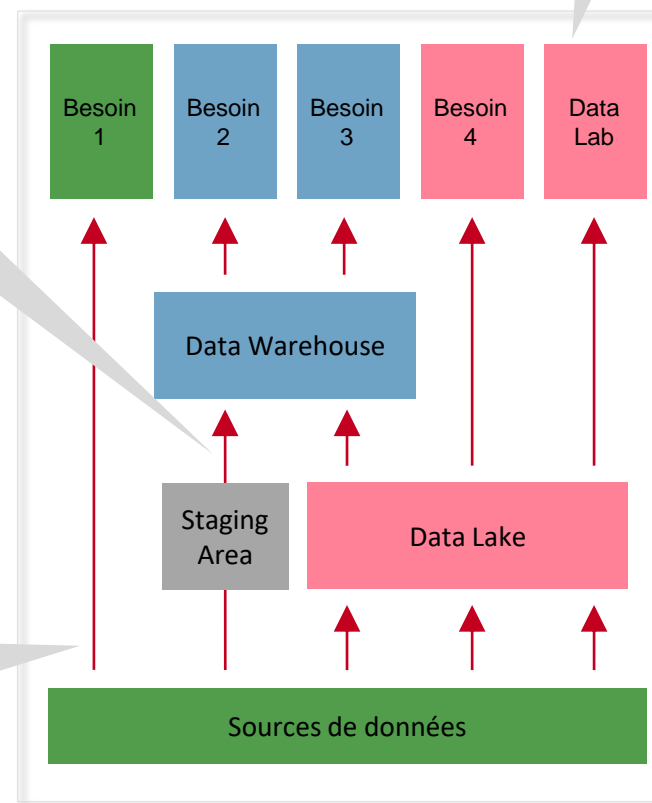
Architecture existante



Les flux de  
données ne seront  
pas rapatriés  
automatiquement  
vers le Data Lake

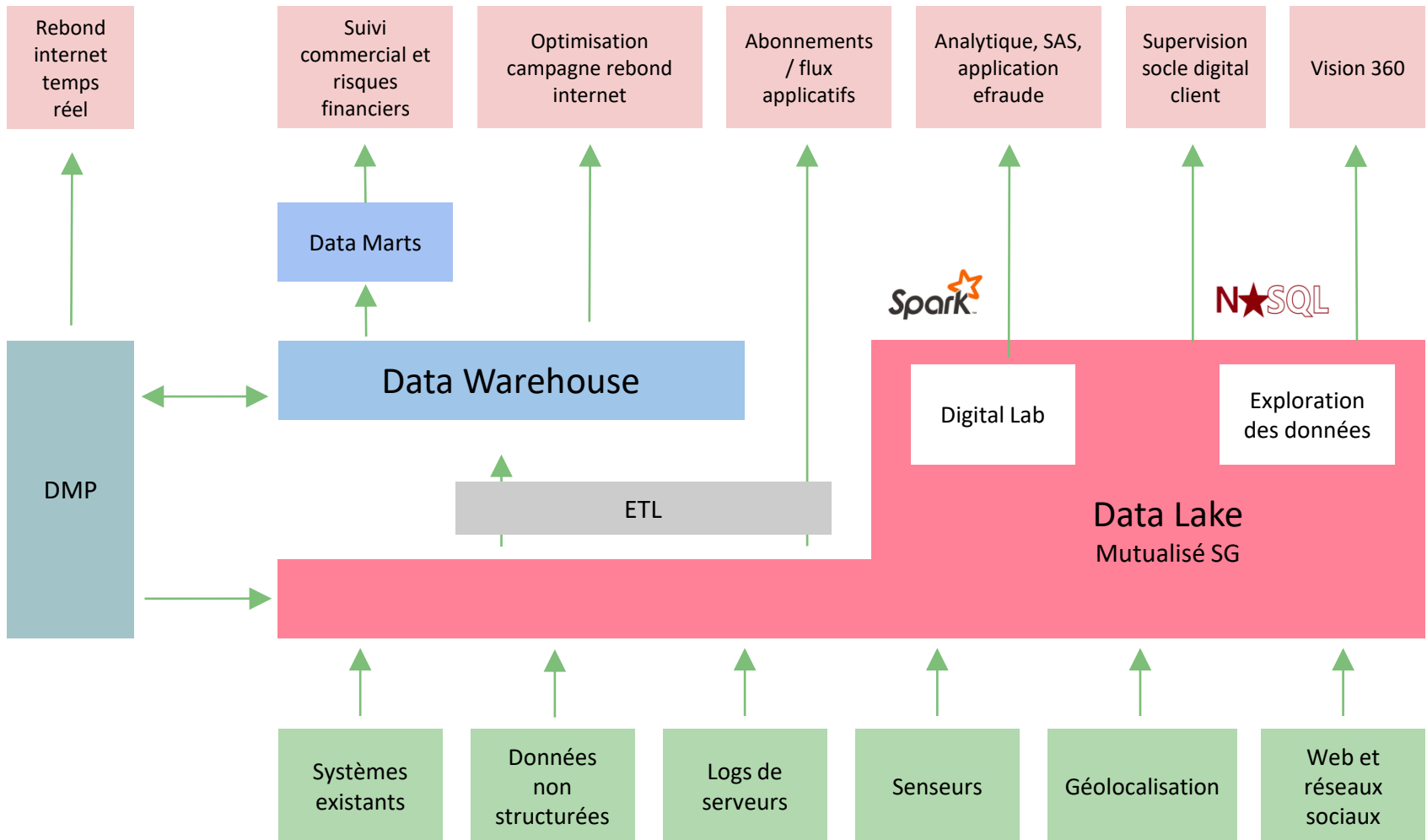
Possibilité de  
donner un accès  
direct aux sources  
(en mode API)

Architecture de référence



## Annexe 4 : CoPil d'architecture Big Data

### Architecture de référence détaillée – Schéma général



- ❑ Ce comité a pour objectif de présenter pour validation les conclusions et les préconisations de l'étude d'architecture sur l'avenir du distributeur de données 3D, réalisée en collaboration avec le BSC .
- ❑ En synthèse :
  - S'appuyer sur 3D est **viaable à court terme**, mais les **risques à moyen terme** sont à considérer :
    - 3D n'a pas la possibilité d'évoluer pour contribuer à la transformation de notre SI (mise à disposition des données pour des usages fil de l'eau / temps réel, mode service, time-to-market, continuous delivery, ...)
    - La technologie sous-jacente à 3D (ETI-Extract) est en obsolescence structurante
  - Devant ce constat, **3 options sont envisageables** :
    1. Ne rien faire
    2. Lancer un projet de migration de 3D en 2017
    3. Mettre sous contrôle l'usage de 3D, limiter son usage et poser les premières briques de la solution alternative à 3D
  - L'étude menée en Q2 2016 a abouti à **préconiser l'option 3**, qui se traduit par
    - Une étude de faisabilité à lancer sur la/les solution(s) d'échanges (adossée à l'écosystème Big Data) qui se substituerait à 3D
    - La mise en place des KPI permettant de suivre l'évolution de 3D au sein d'ITIM
    - Des règles d'encadrement de l'usage de 3D à mettre en place

# Annexe 5 : CoPil d'architecture Avenir 3D

## Architecture d'échange de données : cible

