

# Big Data et données confidentielles

LE CASD ET LA PLATEFORME TERALAB

Alexandre Marty [ [alexandre.marty@casd.eu](mailto:alexandre.marty@casd.eu) ]

# Plan

2

- ▶ Le CASD
  - ▶ Démo
- ▶ La plateforme TeraLab
- ▶ Le projet Données de Caisse
  - ▶ Démo

# Le CASD

# Le Centre d'Accès Sécurisé aux Données

- ▶ Accès et traitement de données confidentielles
- ▶ Activités :
  - ▶ Premier hébergeur français de données administratives pour la recherche
  - ▶ Fournisseur de services pour le public et le privé
- ▶ Plusieurs aspects :
  - ▶ Hébergement, data center
  - ▶ Terminaux, technologie d'accès
  - ▶ Conseil, support, maintenance

# Quelques chiffres

5

- ▶ Equipe d'une vingtaine de personnes
- ▶ 80 sources de données
- ▶ 200 projets
- ▶ 800 chercheurs
- ▶ Projets d'entreprises privées
  
- ▶ Croissance rapide

# Technologie

- ▶ Equipement permettant d'assurer la sécurité de données très sensibles
  - ▶ Authentification forte par biométrie
  - ▶ Prévient physiquement et logiquement toute évacion de fichier de données
- ▶ Equipex
- ▶ Des travaux de valorisation de la technologie

# La SD-Box

7



# Un boîtier sécurisé conçu par nos soins

8

Hautement sécurisée et dédiée à l'accès sécurisé aux données sensibles

Authentification forte et incontournable

Verrouillée physiquement et logiquement, pour limiter la récupération de fichiers



Des boîtiers standardisés

- Fiabilité (configuration unique, stable et validée)
- Déploiement simple et économique
- Besoin très limité d'assistance

Facile à déployer

- Nécessite un écran, un clavier & une souris
- Nécessite une connexion à internet
- Simple à configurer
- Pas d'incidence sur le reste du SI

Pas de données sensibles stockées sur la SD-Box

# La bulle : notre infrastructure étanche

9

La Bulle est **un ensemble étanche de serveurs sécurisés**

Les applications et les traitements de l'utilisateur ne s'exécutent que dans la Bulle

Les insertions / extractions de données sont **contrôlées**. Les utilisateurs n'ont pas accès à internet depuis leur espace de travail.

Les données sensibles sont hébergées uniquement dans la Bulle

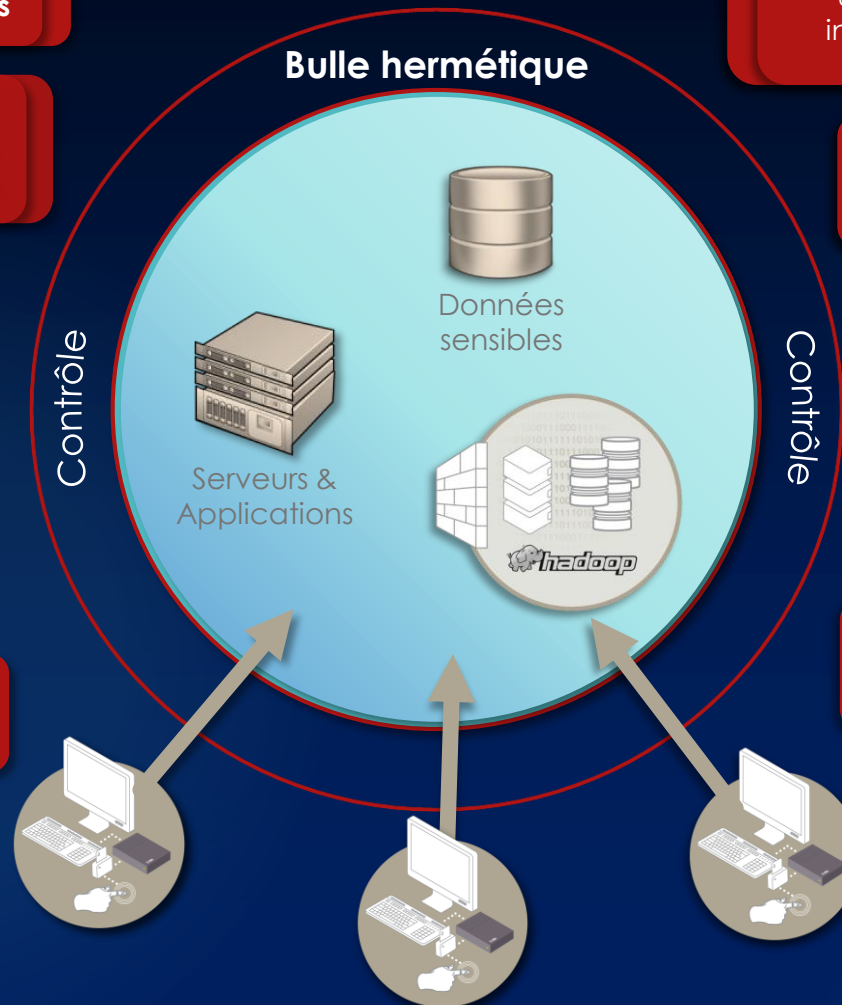
Insertions

Extractions

Les SD-Box sont l'unique moyen d'accéder à la Bulle

Cet accès s'effectue via internet par **canal chiffré**

Un cluster Hadoop est à disposition pour les traitements **BigData**



# La SD-Box

10



← La bulle

← La SD-Box

← Les données

← Les outils

# L'offre classique CASD

11

- ▶ Tous les outils dont les chercheurs ont besoin
  - ▶ Environnement Windows
  - ▶ SAS, R, Stata, SPSS, Python...
  - ▶ Ajout de tout logiciel sous réserve de licence compatible
- ▶ Plug & Play
  - ▶ Il suffit de brancher la SD-Box et ça marche
  - ▶ Support technique

# Démo

# Plateforme TERALAB



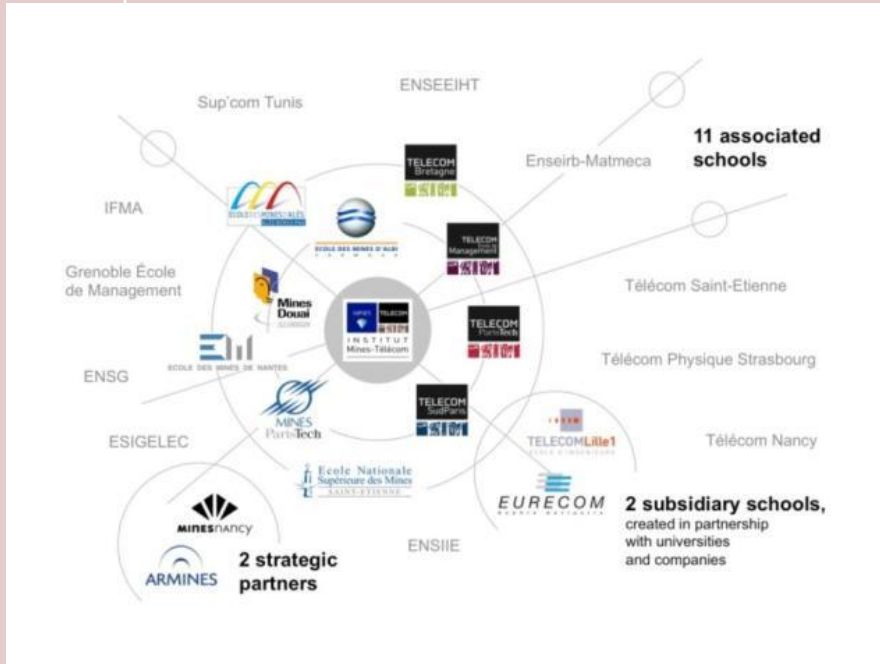
# Le projet TERALAB

14

- ▶ Appel à projets “Cloud computing/Big Data” dans le cadre du PIA (Programme d'Investissements d'Avenir)
- ▶ Construction et exploitation d'une plateforme Big Data
  - ▶ Pour la recherche, l'innovation et l'enseignement
  - ▶ Soumise par un consortium comprenant
    - ▶ L'[IMT](#) (Institut Mines-Télécom)
    - ▶ Le [GENES](#), et particulièrement le CASD
    - ▶ En partenariat avec l'[INSEE](#)
- ▶ Projet sélectionné et lancé
  - ▶ Budget de 5,7 M€
  - ▶ Durée de 5 ans
  - ▶ Signature en Décembre 2013

# Acteurs

## Institut Mines-Télécom



## GENES

 Enseignement - Cycle Ingénieur - Master	 Enseignement - Cycle Ingénieur - Master	 Centre d'Accès Sécurisé aux Données (CASD)
Recherche 7 laboratoires	Recherche 2 laboratoires	 Conseil & Expertise
Formation professionnelle	Formation professionnelle	 Institut des Politiques Publiques Institut des Politiques Publiques (Ecole d'Économie de Paris et GENES)
 ÉCOLE NATIONALE DE LA STATISTIQUE ET DE L'ADMINISTRATION ÉCONOMIQUE - PARIS Campus Paris-Saclay	 ÉCOLE NATIONALE DE LA STATISTIQUE ET DE L'ANALYSE DE L'INFORMATION - GENNES Campus de l'Université Européenne de Bretagne	 Coopération internationale et Appui aux Ecoles de statistique étrangères

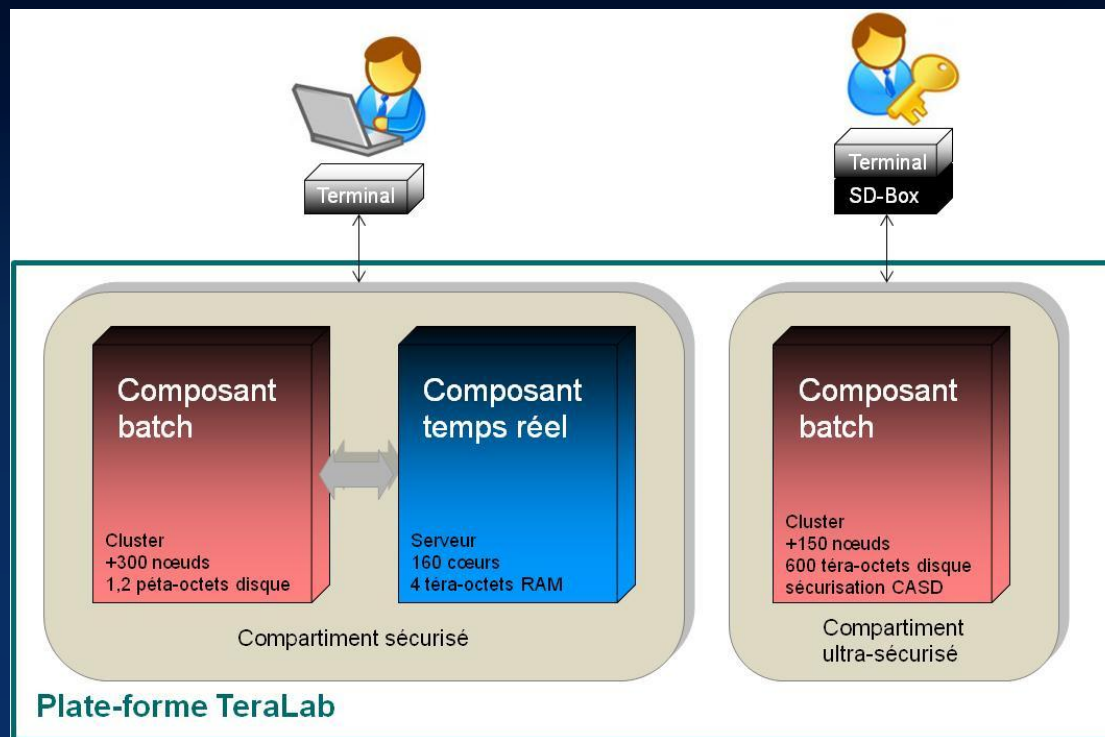
**Domaines majeurs** Numérique, énergie, environnement  
Matériaux, économie/ entreprise/ société

**Tutelles** Ministère du redressement Productif, de l'Enseignement Supérieur et de la recherche

Economie, Sciences humaines et sociales  
mathématiques appliquées, Statistiques, économétrie

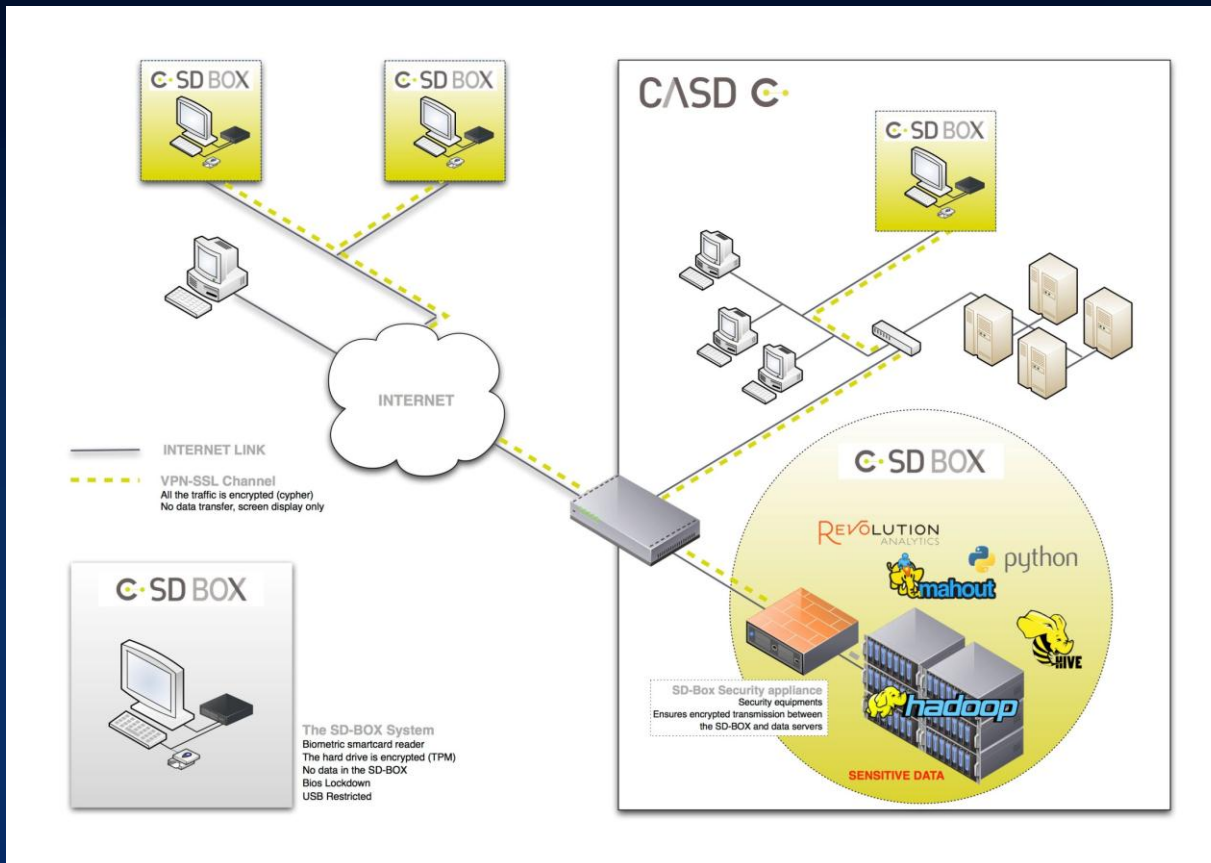
Ministère de l'Économie et des finances (INSEE)

# Compartiments de la plateforme



# Architecture TERALAB-CASD

17



# Technologies Big Data

18

- ▶ Infrastructure puissante et extensible
- ▶ Environnement Hadoop
- ▶ Nombreux outils pour les data scientists :
  - ▶ Python, R
  - ▶ Machine learning
  - ▶ Serveur Open Street Map

- ▶ Données de caisse
- ▶ RTE (Réseau de Transport d'Electricité)
  - ▶ Problématiques d'optimisation du réseau (prévision, maintenance et consommation)
  - ▶ Grande diversité de sources
  - ▶ Développement d'applications innovantes
- ▶ En discussion : Données de santé
  - ▶ Sujet en pleine actualité
  - ▶ Grande confidentialité
- ▶ Implication dans des projets européens : DwB, Eurostat Big Data Task Force

# Projet Données de Caisse

- ▶ Données de caisse de grande taille
- ▶ Pertinence des technologies Big Data
- ▶ Cluster Hadoop ultra sécurisé du CASD
- ▶ Outils scientifiques puissants

- ▶ Amélioration de l'IPC :
  - ▶ capture des prix réels des ventes
  - ▶ exhaustivité
- ▶ Nouvelles productions, en gardant les concepts actuels de l'IPC :
  - ▶ prix moyens
  - ▶ indices régionaux
  - ▶ comparaisons spatiales

- ▶ Fichier par semaine et par enseigne avec les mesures suivantes :
  - ▶ Prix de vente
  - ▶ Quantité
  - ▶ Point de vente
  - ▶ Code barre (EAN)
  - ▶ Date de vente
  - ▶ Famille de l'article
  - ▶ Description de l'article
- ▶ Prix et quantités simulés
- ▶ Référentiels des points de vente + Référentiel des articles : volume faible

# Cluster actuel

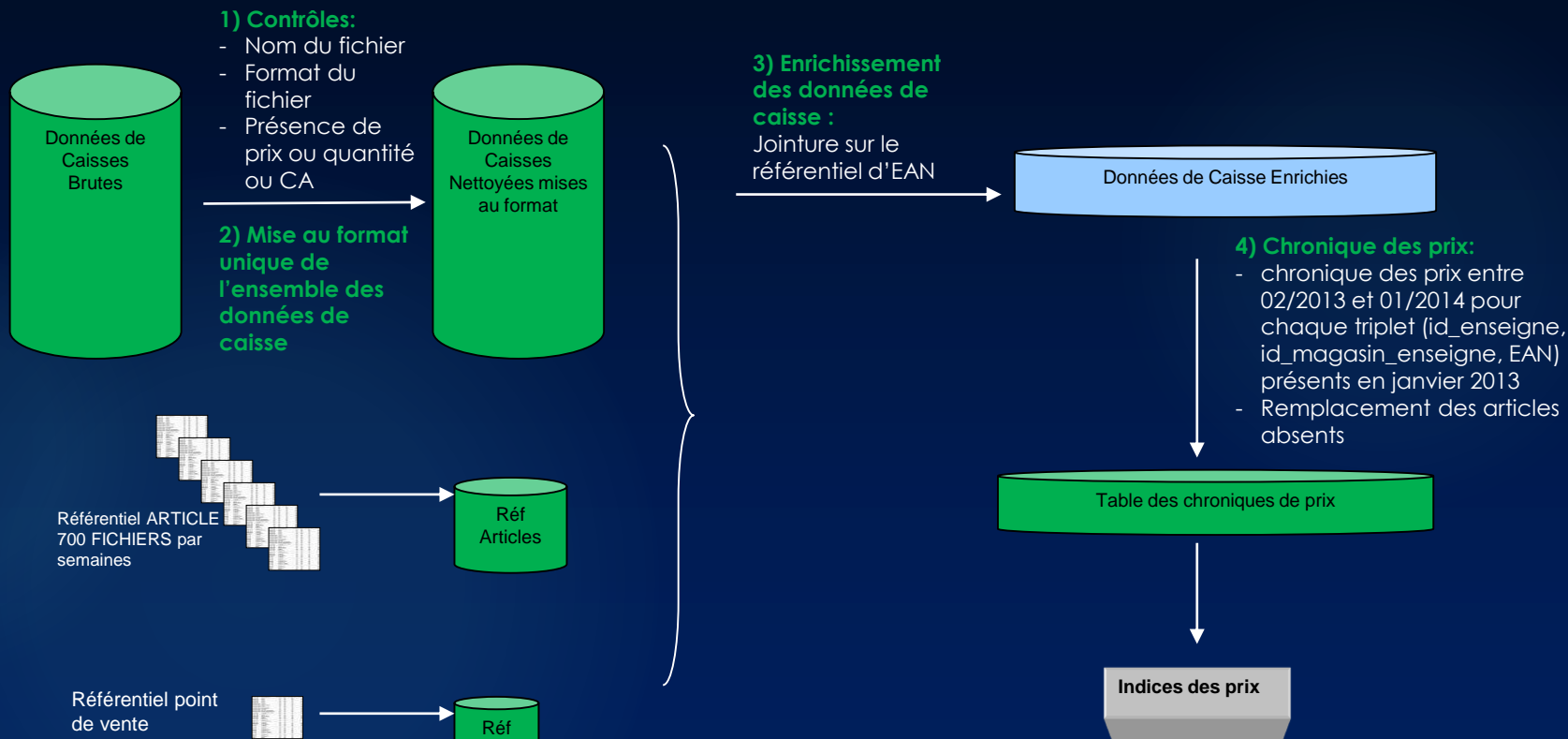
24

- ▶ 6 machines : 2 nœuds de gestion , 4 nœuds de travail
- ▶ 48 cœurs, 512 GB de RAM
- ▶ 33 TB de stockage
  
- ▶ Extensible à la volée

- ▶ Situation réelle actuelle :
  - ▶ 5 enseignes livrent 10 Go de données par semaine
  - ▶ 900 Go de données à traiter
  - ▶ 5,7 milliards de lignes
- ▶ Situation de test :
  - ▶ Situation réelle x5
  - ▶ 28,5 milliards de lignes
- ▶ Situation de long terme : 10x situation actuelle

# Traitements

26



Formule de Laspeyres:

IPC par famille de produits =

$$\prod_{i=1}^n \left( \frac{P_i^t}{P_i^0} \right)^{ca} / \sum ca$$

→ **IPC**

- ▶ Pig, Hive, Impala :
  - ▶ Chargement et traitements type base de données SQL
- ▶ Python :
  - ▶ Orchestration des tâches

- ▶ Temps de traitement :
  - ▶ Chargement de 70 semaines de données : 5h
  - ▶ Constitution du panier : 2h40
  - ▶ Jointure référentiel produit : 3h05
- ▶ Par rapport à solution Oracle actuelle :
  - ▶ Complexité et temps de développement très réduits
  - ▶ Temps de traitement très inférieurs

Merci !

Questions ?

Alexandre Marty

[ [alexandre.marty@casd.eu](mailto:alexandre.marty@casd.eu) ]