



# BIG DATA PLATFORM FOR RESEARCH AND INNOVATION

## TERALAB

Présentation



TERALAB

DATA SCIENCE FOR EUROPE



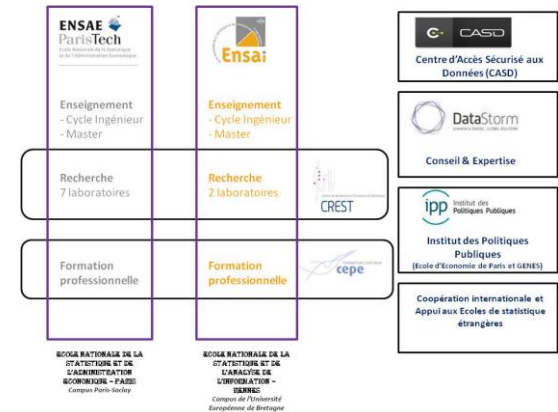
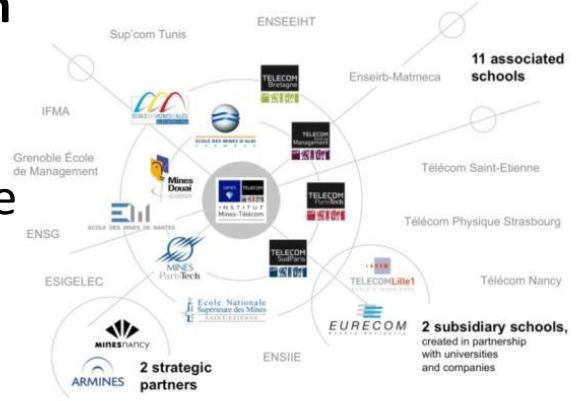
cap·digital



**Teralab**

**Vocation : Encourager l'exploitation massive des données**

- **Une plateforme dédiée à la recherche / innovation / enseignement**
  - projets de recherche – innovation,
  - Enseignement, dont soutien formation continue
  - ... en dehors de toute exploitation commerciale.
- **Dan un cadre**
  - De projet collaboratif
  - Ou bilatéral
- **Financée par le PIA (5,7M€)**
  - 2012 : appel à projet PIA remporté par l'**IMT** et le **GENES**, en partenariat avec l'INSEE
  - Fin 2013 : validation du volet financier par la caisses des dépôts et consignation
  - Pour une durée de 5 ans (T0 janvier 2014)
- **La plateforme a été ouverte en avril 2014**





**Teralab**

## **Un positionnement de catalyseur**

- Teralab se positionne en catalyseur / accélérateur de l'exploitation des données par les acteurs économiques

**1. Lève les freins pour les acteurs hésitants à passer à l'action sur leurs données**

**2. Accélère la maturation des usages grâce à l'orientation innovation / R&D comme cadre imposé**

**3. Transfère les projets matures aux acteurs concurrentiels dès lors que le partenaire sait ce dont il à besoin**

**4. Promeut les acteurs innovants de l'écosystème : installés sur la plateforme pour faciliter la rencontre avec leur marché**



**TERALAB**

DATA SCIENCE FOR EUROPE



**cap-digital**  
Paris Region



## Teralab / cadre d'intervention

### Synthèse des modalités d'accès

#### DISPOSITIFS \*

##### PROJETS R&D COLLABORATIFS

- Des consortiums R&D composés d'industriels, de laboratoire de recherche et de PME/startups

##### CHALLENGES BIG DATA CAP DIGITAL

- Des grands groupes souhaitant être sponsors de l'initiative
- Des PME/Startups big data souhaitant avancer avec un sponsor fournisseur de données

##### POCs BILATERAUX

- Des consortiums commerciaux composés de grands groupes, cabinets de conseils et startups big data

##### PROJETS "BOOTSTRAP"

- Un consortium R&D ou commercial intégrant un chercheur de l'institut Mines Télécom

##### ENSEIGNEMENT

- Des projets type fil rouge académiques ou industriels souhaitant proposer une alternative souveraine à un hébergement classique type aws

#### CIBLES

#### CONDITIONS DE FINANCEMENT

- label R&D requis, national ou européen (ANR, FUI, H2020, Carnot...)
- Sélection par le comité scientifique TeraLab
- Financement variable (jusqu'à 100K€ environ)

- sélection des challenges par le consortium
- Equipe projet est co-financée par l'Etat et le Sponsor à hauteur de 100K€ par an en moyenne (sur une durée de 3 ans)

- Sélection par le comité scientifique TeraLab
- Pas de financement direct mais possibilité d'inclure d'autres financements (thèses CIFFRE, labos communs)

- Le projet est financé à hauteur de 100K€ maximum par l'intermédiaire du laboratoire de recherche
- Sélection par le comité scientifique TeraLab
- La décision finale revient à la DGE

- Sélection par le comité scientifique TeraLab

\*document de travail du 06/11/2014

TERALAB

DATA SCIENCE FOR EUROPE



cap.digital



## Teralab / composantes de l'offre

### Panorama global

#### Infrastructure technique

- Hébergement souverain
- compartiments sécurisés et dédiés
- Double environnement : distribués et in memory

#### Accompagnement

- Technique prise en main
  - Scientifique
  - Ecosystème

#### Ecosystème embarqué

- Les solutions standards du marché,
- Les solutions de start up innovantes

#### Corpus de données à disposition

- Etalab
- Crawling ..
- ....

Cible mi 2015



TERALAB

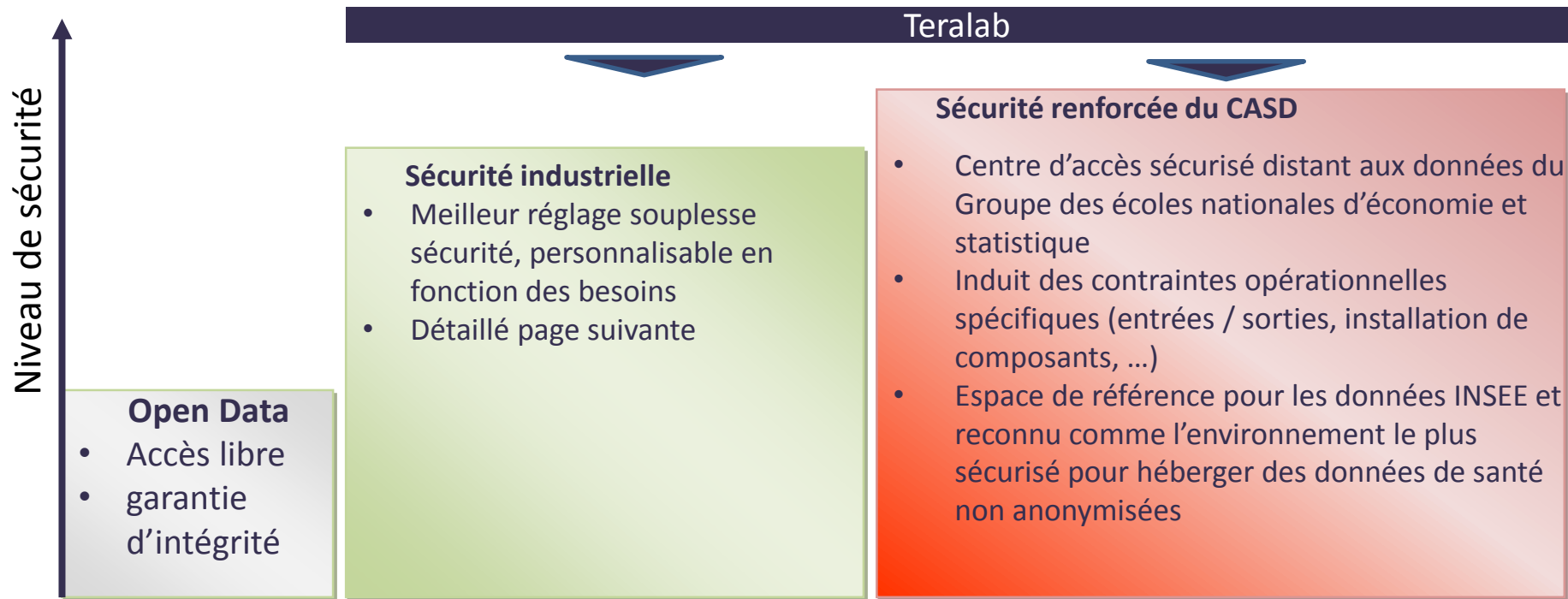
DATA SCIENCE FOR EUROPE



cap digital



- L'activité de recherche / innovation induit une acceptation juridique du risque, cependant Teralab met en œuvre tous les niveaux de sécurité industrielle, déclinée en deux niveaux :



**Une garantie d'hébergement physique en France métropolitaine**

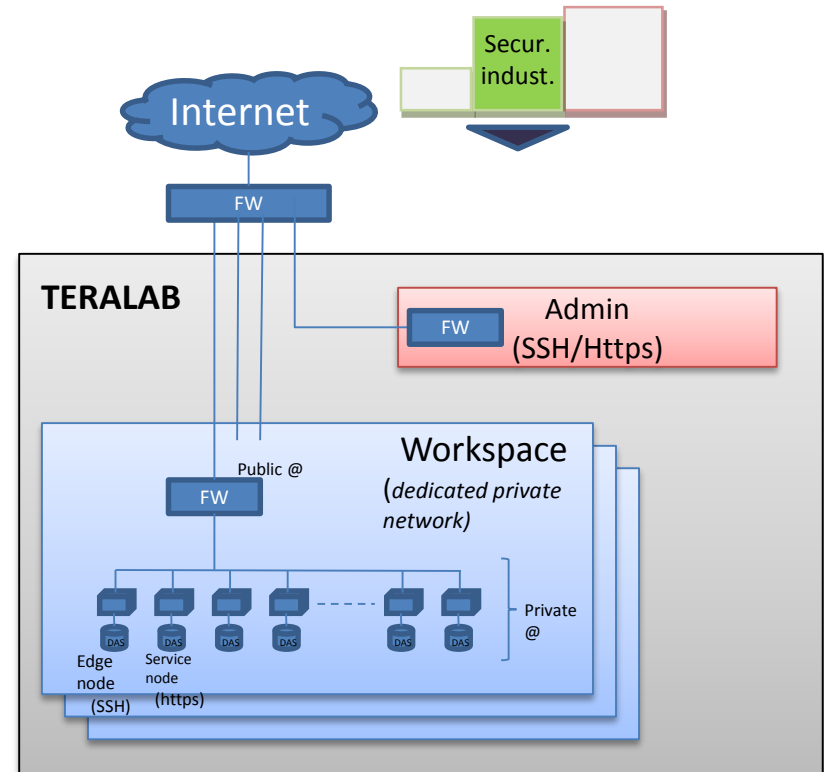




## Le niveau « sécurité industrielle » optimise le compromis accessibilité / sécurité



- Espaces de travail cloisonnés (VLAN)
  - Switch réseau configurés pour constituer des VLAN isolés pour chaque cluster
  - Protège des incursions / perturbations extérieures
- 2 niveaux de firewall
  - En amont de la plateforme : le seul à avoir une IP publique
  - En amont de chaque cluster, personnalisable (sur demande)
- Tout accès nécessite un compte utilisateur (login / mot de passe)
- Droit d'administrateur pour le partenaire lui permettant de personnaliser son dispositif de protection (authentification par clé par exemple)



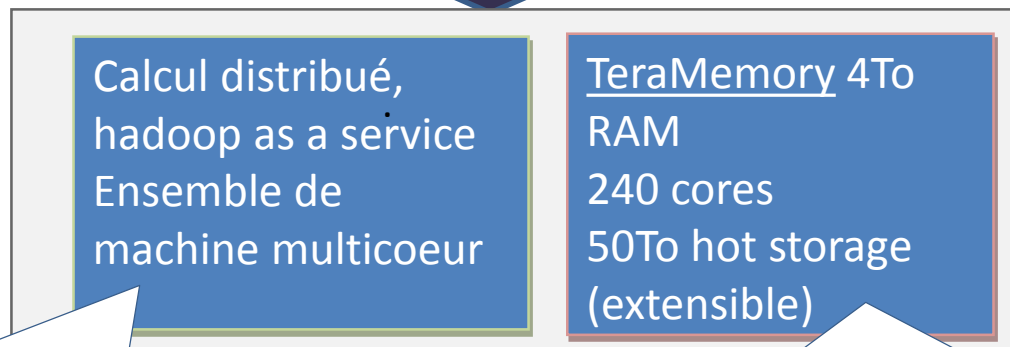
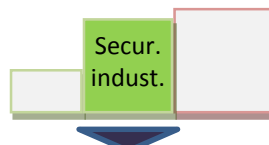
- Le niveau de sécurité peut être ajusté en fonction des besoins d'échange de données du partenaire (objets connectés par exemple)





## Teralab / offre - infrastructure

### 2 infrastructures de calcul complémentaires



#### exemples de configuration de cluster en production

Hadoop	vcpu	RAM (GB)	Storage (GB)
CDH5.0	21	148	3512
CDH5.0	13	68	1944
CDH5.0	69	404	14744
HW	13	68	3512

Data provisionning  
SQL / No sql / hadoop

#### Configurations à très fort ratio mémoire / cœur

- Memory upgradable from 4 to 24 TeraByte
- 240 cores (8 Modules x2 Xeon/module x15cores/Xeon)
- 8000 specint
- Net storage (DAS - Raid5):
  - current TeraLab: 50TB /
  - Max : 800TB



8-9 octobre 2014



TERALAB

DATA SCIENCE FOR EUROPE



cap.digital





Teralab / offre - écosystème embarqué

Des ressources installées / une personnalisation par vos soins



- Les deux compartiments e calcul sont prééquipés de l'infrastructure logicielle au niveau de l'état de l'art

**Linux** : CentOS 6.4

**Hadoop** : Cloudera Full CDH5 ou HW2 (dont web Hue, ClouderaMgr,)

**Python 2.7** : Canopy (Scikit-learn installable), **R**, **Dataiku**, **Pentaho** (partie libre)

**Spark, Flink** (ex stratosphère)

**VNC**

**Ganglia**

Stockage et  
calcul distribué

analytique

Calcul distribué  
in memory

Bureau déporté

Monitoring

Les briques installées sont

- open source
- briques innovantes de start up (ex dataiku) a des fins de promotion

Les gestionnaires de projets hébergés ont un compte administrateur et peuvent **installer toute brique logicielle pour lesquelles ils disposent d'une licence**



TERALAB

DATA SCIENCE FOR EUROPE



cap.digital



## Teralab / offre - écosystème embarqué

### Exemple d'une configuration Teralab

#### WORKSPACE (project, team...)

##### Users

- Roles (applis, admin. actions)
- Groups (files)



##### Edge nodes

- Linux server
- unix



##### Cluster

- Create, start, stop, resize, opera
- HDFS, unix



**Edge node:** The Edge Node is a system intended to be used as a Gateway. It is not a real part of the Cluster but all client library are installed on, thus allowing full access to the cluster from this system.



**Service node:** System hosting most of helper services (Hue, Ganglia, Cloudera Manager, Ambari ...).



**Name nodes:** Two Namenodes acting as active/passive High Availability configuration.



**Data nodes:** The workhorses of the cluster. Provide storage and compute resources.



**Firewall:** A pfSense Firewall controlling all cluster access.

#### Environment

##### Hadoop Environment

- Hadoop: Cloudera 4.6 + Cloudera 5

##### Administration

- TeraLab mgr (cluster), Nagios (alerts), Ganglia (monitoring)

##### User operation

- Hue, ssh

##### Applications

- Oozie, Pig, Hive, Impala, Mahout...
- Dataiku, Opendatasoft...

##### Database

- HBase (columnar)

##### Parallelisation Engine

- Yarn, MapReduce, Stratosphere
- Spark

##### File system

- HDFS

##### Cluster OS

- CentOS 6.4 (incl. python 2.6.6)
- Others linux distributions

##### Import/Export

- Flume, Hive Metastore Mgr, Sqoop

##### Others

- Any imported Linux tool (may run on the cluster or on its edge node)



TERALAB

DATA SCIENCE FOR EUROPE



cap.digital



Teralab / offre - Accompagnement

Une équipe de 6 personnes dédiée en soutien  
des utilisateurs de la plateforme



- **Sur le volet technique**
  - Assistance à la prise en main de l'infrastructure : conseils et aide opérationnelle
  - Hotline avec des tickets 'incidents en ligne
- **Sur le volet scientifique**
  - Une directrice scientifique dédiée au projet, avec une très forte expérience opérationnelle
  - Un lien direct avec les équipes de recherche Mines et Telecom (dont laboratoires de machine learning
- **Sur le volet économique**
  - Une orientation si besoin au sein de l'écosystème économique Big Data : acteurs du marché, financements potentiels, start up sur des activités connexes



## Plateforme référencée dans le cadre du plan Big Data (Plan “nouvelle France Industrielle”)

..





# Teralab – écosystème

## Projets déjà hébergés

### Project description

### Actors involved

#### Pure R&D

- Machine Learning project (?)
- SickitLearn project (?)

- Big Data Researcher: A.Gramfort, S.Clemencon, ...



#### Collaborative R&D

- Speedata - Decision-making platform in SaaS mode
- D4D 'Data for Development Senegal' - innovation challenge open on ICT Big Data for the purposes of societal development in Senegal.
- Square Predict - Big Data SaaS platform for insurance industry.
- Normatis - Use of multimodal data with Web information extraction, data integration, and mobility data management.

- Large industrial groups, SMEs, R&D Labs:



#### Teaching

- Yearly project (?)
- Lab work project (?)
- Others (?)

- Lecturer and researcher, students of big data masters:



#### Industrial Proof of Concept

- Proof of Concept– Fraud Detection for the french actor of private health insurance
- Proof of Concept Pôle Emploi (French Employment Office) – Matching of jobs offers/profile/specific training to improve unemployed people market reinsertion

- Large industrial groups with business/big data challenges:





## 4 chaires dans le périmètre Mines Telecom directement en lien avec le Big Data

	Pilote	Partenaires
Machine Learning for Big Data	Stéphan Cléménçon	BNP Paribas , Criteo , PSA, Safran
Valeur et politique des informations personnelles	Claire Levallois Barthes	Groupe Imprimerie Nationale, BNP Paribas et Dassault Système, avec l'aide de la CNIL
réseaux sociaux	Christine Balagué	la Poste, Pages Jaunes Groupe et Danone
E commerce	Talel Abdessalem	Partenaires



TERALAB

DATA SCIENCE FOR EUROPE



cap·digital