



# Tour de table

## Tour de table : 3 questions

- Présentation personnelle Qui ? (Je suis ...)
- Entreprise/Contexte professionnel
- Attentes (les attentes de chacun par rapport à la formation)







# Sommaire

1. **Généralités**
2. Les opportunités qu'offre le Big Data pour la DSI
3. Open Data
4. Cas d'usage du Big Data
5. Le stockage dans le Big Data
6. Les technologies du Big Data
7. Le traitement des données en Big data
8. Compétences autour du Big Data
9. Les étapes d'un projet Big Data



Explosion de la quantité des données

Le partage des données

La recherche des données

Le stockage des données

Le traitement des flux de données

# Le partage des données

# La recherche des données

# Le stockage des données

# Le traitement des flux de données





# Qu'est-ce que Big Data ?

Big Data désigne un ensemble très volumineux de données qu'aucun outil classique de gestion de base de données ou de gestion de l'information ne peut manipuler convenablement.

Exemples :

- ✓ environ 2,5 trillions d'octets de données chaque jour.
- ✓ ce sont des informations provenant de sources diverses :
  - messages
  - vidéos
  - informations climatiques
  - signaux GPS
  - transactions
  - etc.





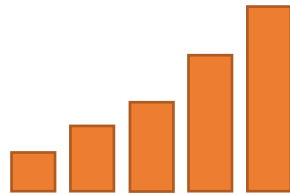


# Qu'est-ce que Big Data ?

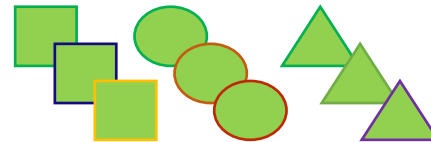
- La taille n'est pas l'unique caractéristique.
- Flux de données continu (big data stream).
- Plus personne ne veut jeter les données.
- La variété des sources de données indépendantes nécessite l'intégration des données.
- Hétérogénéité des données.
- Souvent associé aux données non structurées, mais non exclusivement.



# Volume



## Variété



## Vélocité





# Qu'est-ce que Big Data ?

Plus **2 V**

**V**éracité

**V**aleur

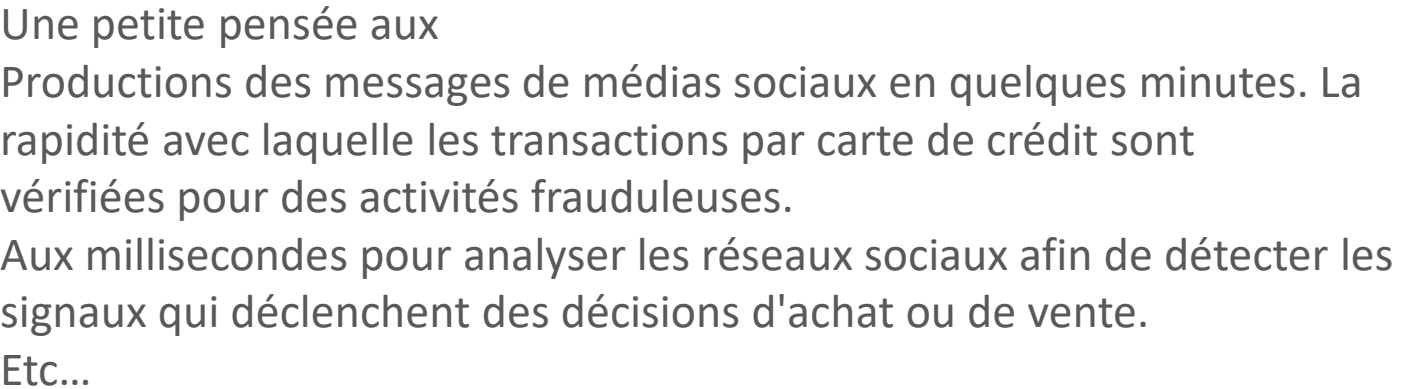






La vitesse fait référence à la vitesse

- ✓ de la génération des nouvelles données
- ✓ à laquelle les données se déplacent

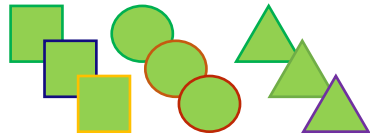


La technologie Big Data nous permet aujourd'hui d'analyser les données pendant qu'elles sont générées sans les stocker dans des bases de données.



# Qu'est-ce que Big Data ?

## Variété



La variété fait référence aux différents types de données utiliser.

Avant :

Stockage des données structurées dans les bases de données relationnelles.

Aujourd'hui :

80% des données ne sont pas structurées

Difficile de les stocker dans des bases de données relationnelles

- ✓ photos
- ✓ vidéo
- ✓ audio
- ✓ messages réseaux sociaux

Avec la technologie Big Data nous pouvons désormais exploiter différents types de données, et les rassembler avec des données structurées et traditionnelles.



# Qu'est-ce que Big Data ?

## Véracité



La véracité fait référence au désordre ou à la fiabilité des données.

La qualité et la précision sont moins contrôlables

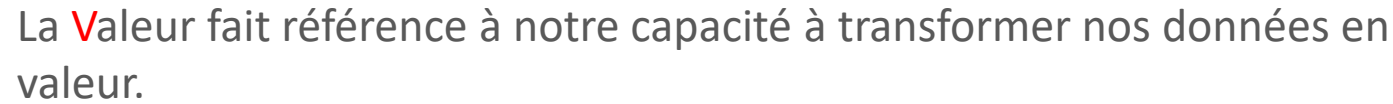
- ✓ messages Twitter avec des hashtags
- ✓ les abréviations
- ✓ les fautes de frappe
- ✓ Etc...

La technologie Big Data et Analytics nous permet désormais de travailler avec ces types de données.





**Valeur**



- ✓ Volume
- ✓ Variété
- ✓ Vitesse
- ✓ Vérité

doivent être transformés en Valeur !

La **V**aleur est le **V** du Big data qui compte le plus.



# Pourquoi Big Data ?

# Difficulté de traiter la masse de données produite chaque jour

# Les informations sont produites en temps réel

# Problème d'optimisation de base de données

# Le traitement des données non structurées



19



L'ensemble du marché mondial des logiciels devrait générer des revenus de 628 milliards de dollars, dont 302 milliards de dollars provenant des applications.





Les initiatives Big Data axées sur ce domaine ont également le taux de réussite le plus élevé (69%) selon la dernière enquête NewVantage Venture Partners.

Plus du tiers des entreprises, soit 36%, affirment que cette zone est leur priorité absolue en matière d'analyse avancée et d'investissement dans le Big Data.



# Freins de la mise en place du Big data

Coût,

# Manque de compétences,

Manque de visibilité sur les opportunités,

# Difficile de quantifier le ROI des investissements Big Data,

La collecte des données doit surpasser les canaux traditionnels,

Les données sont non structurées (nécessité de compétences pour les traiter).



## Adoption de nouvelles technologies d'exploitation de la donnée.





## Les acteurs historiques de solutions IT :

# Oracle

# SAP

25



# Les acteurs du Big Data

## **IBM :**

- ✓ BigInsights Enterprise

## **Microsoft :**

- ✓ privilégié l'utilisation du framework Hadoop
- ✓ Windows Azure et Window Server

## **Amazon :**

- ✓ Amazon Web Services
- ✓ Elastic MapReduce (EMR)





On retrouve également les spécialistes de l'analytique :

- ou encore les fournisseurs spécialisés tels que :

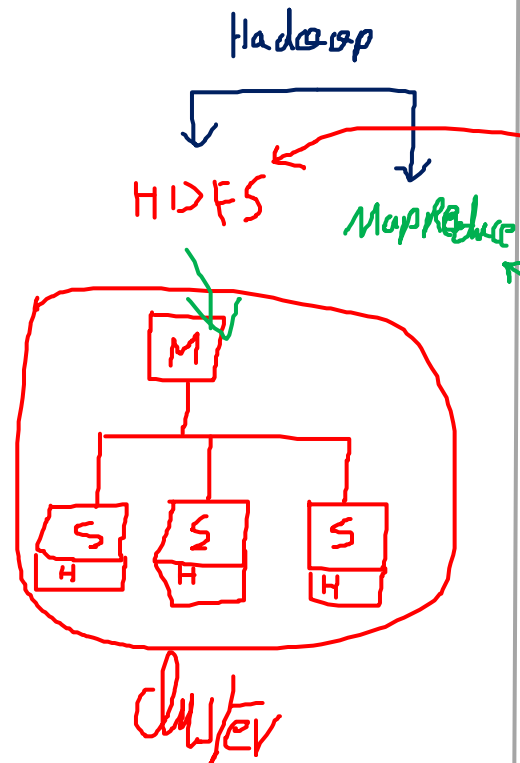
- 28



# Distributions majeures

Hadoop est le Framework logiciel open source au cœur de la révolution Big Data.

- ✓ sortie en 2011
- ✓ solutions pour le stockage
- ✓ l'analyse de données
- ✓ Avec distributions commerciales







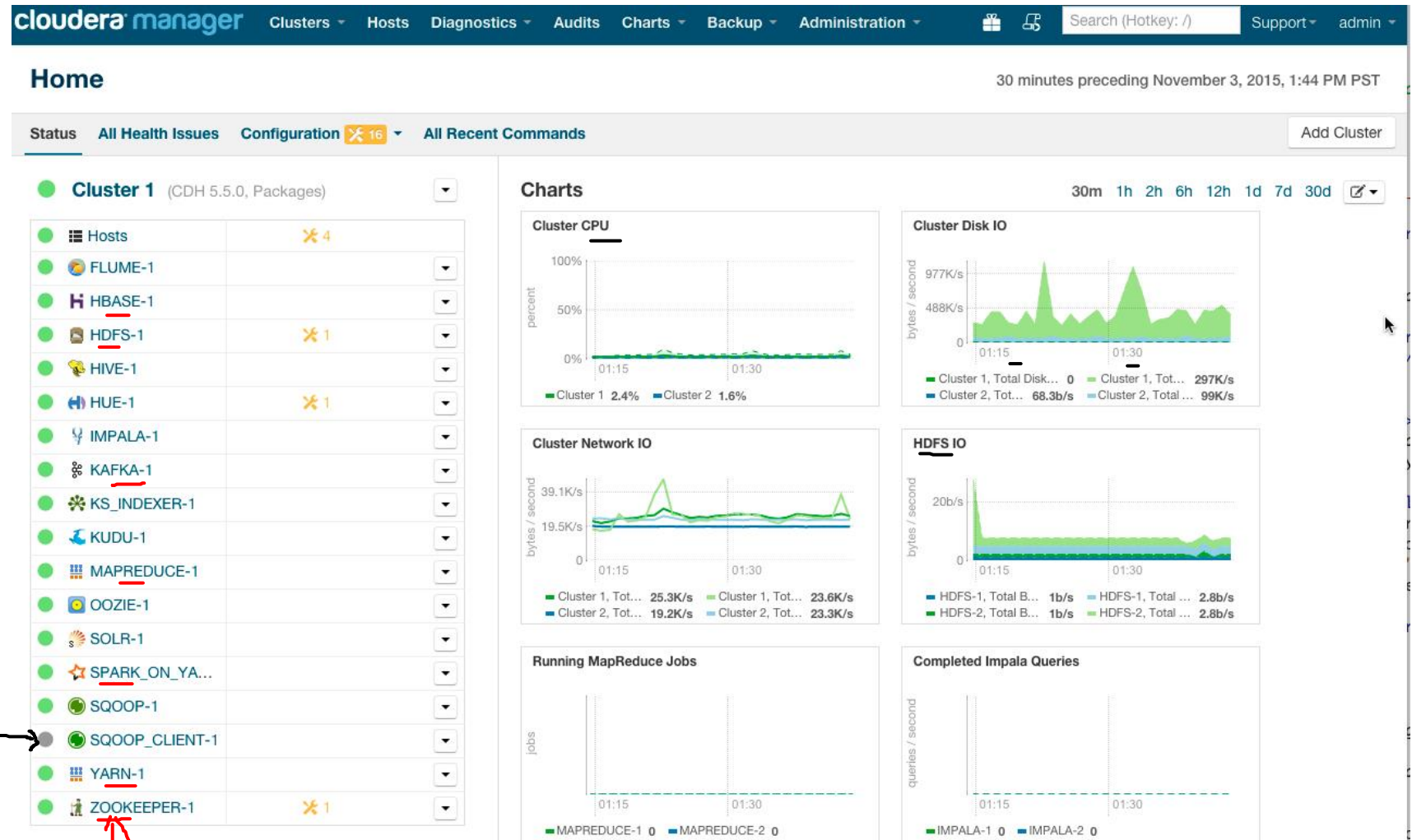
Cloudera a été le premier fournisseur à proposer Hadoop en tant que package et continue d'être un leader dans l'industrie.

# Cloudera CDH avec des composants open source

## La sécurité et les interfaces pour l'intégration avec des applications tierces



# Distributions majeures : Cloudera



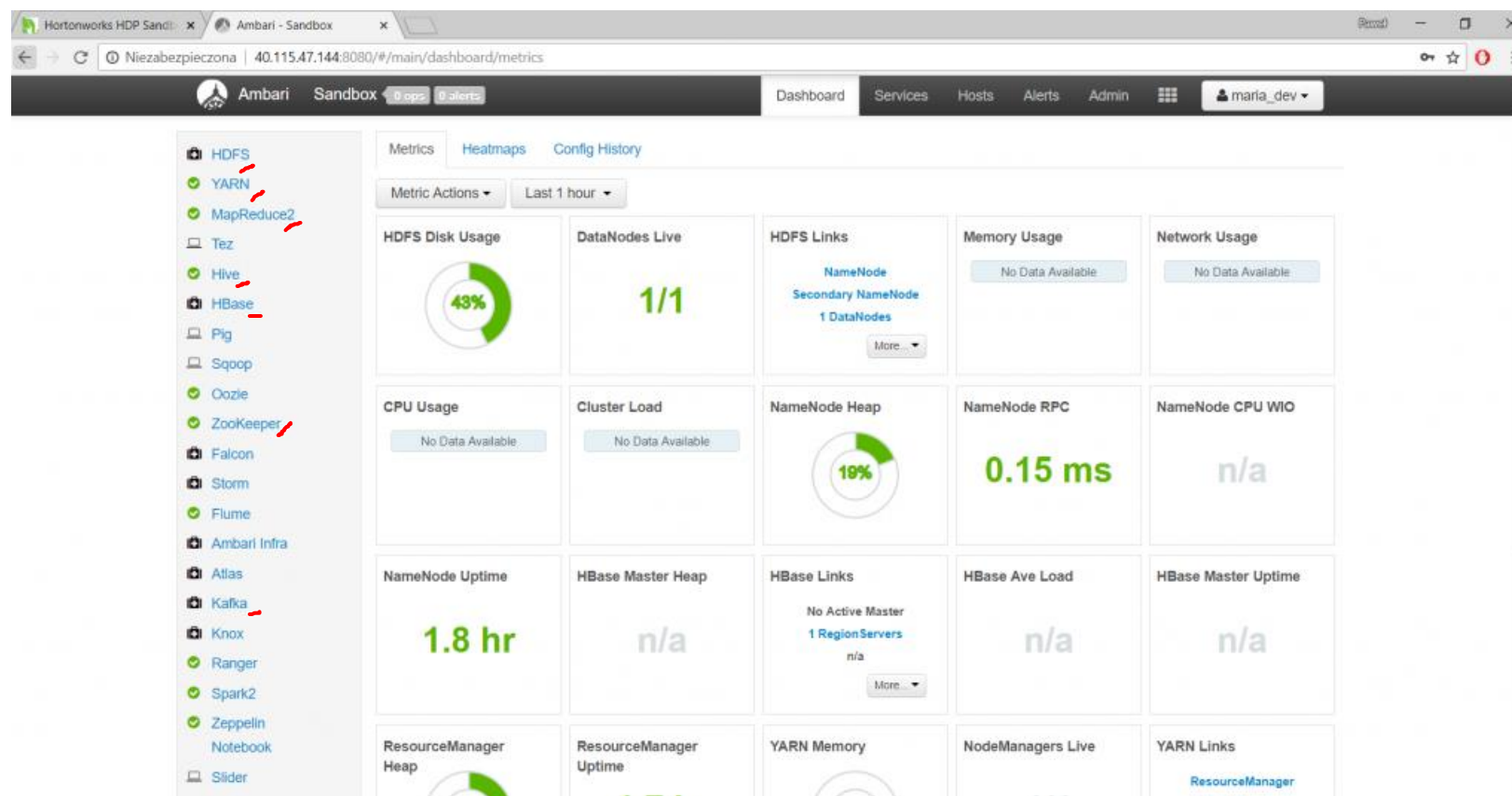




# Amazon et IBM proposent désormais Hortonworks en tant qu'options sur leurs propres plateformes



# Distributions majeures : Hortonworks







36

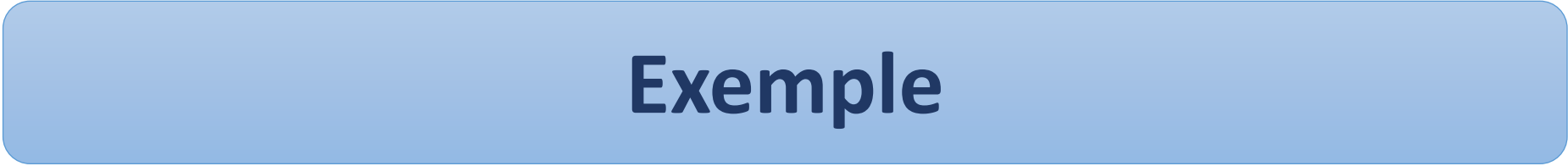


# Distributions majeures : Microsoft

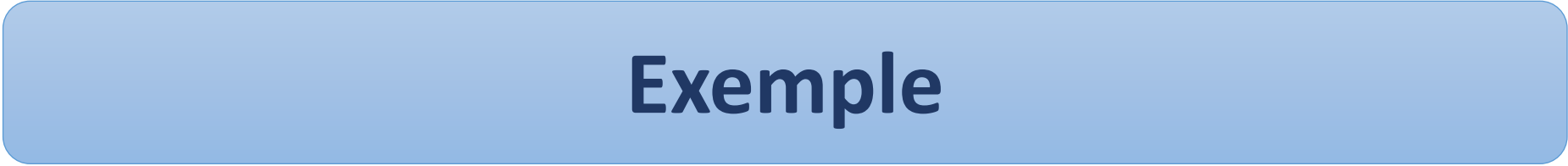
- Azure HDInsight
  - ✓ service en cloud
  - ✓ propose des installations gérées de plusieurs distributions Hadoop (Hortonworks, Cloudera et MapR)
  - ✓ Il les intègre à sa propre plate-forme Azure Data Lake pour offrir une solution complète de stockage et d'analyse basée sur le cloud.
  - ✓ Fournit les services cloud Spark, Hive, Kafka et Storm











# Pourquoi utiliser le Big Data ?

## Kilobyte (KB)

1 Kilobyte = 1024 Bytes

1 KB = Paragraph

100 KB = Low-resolution photo

128 KB = Memory of a standard calculator

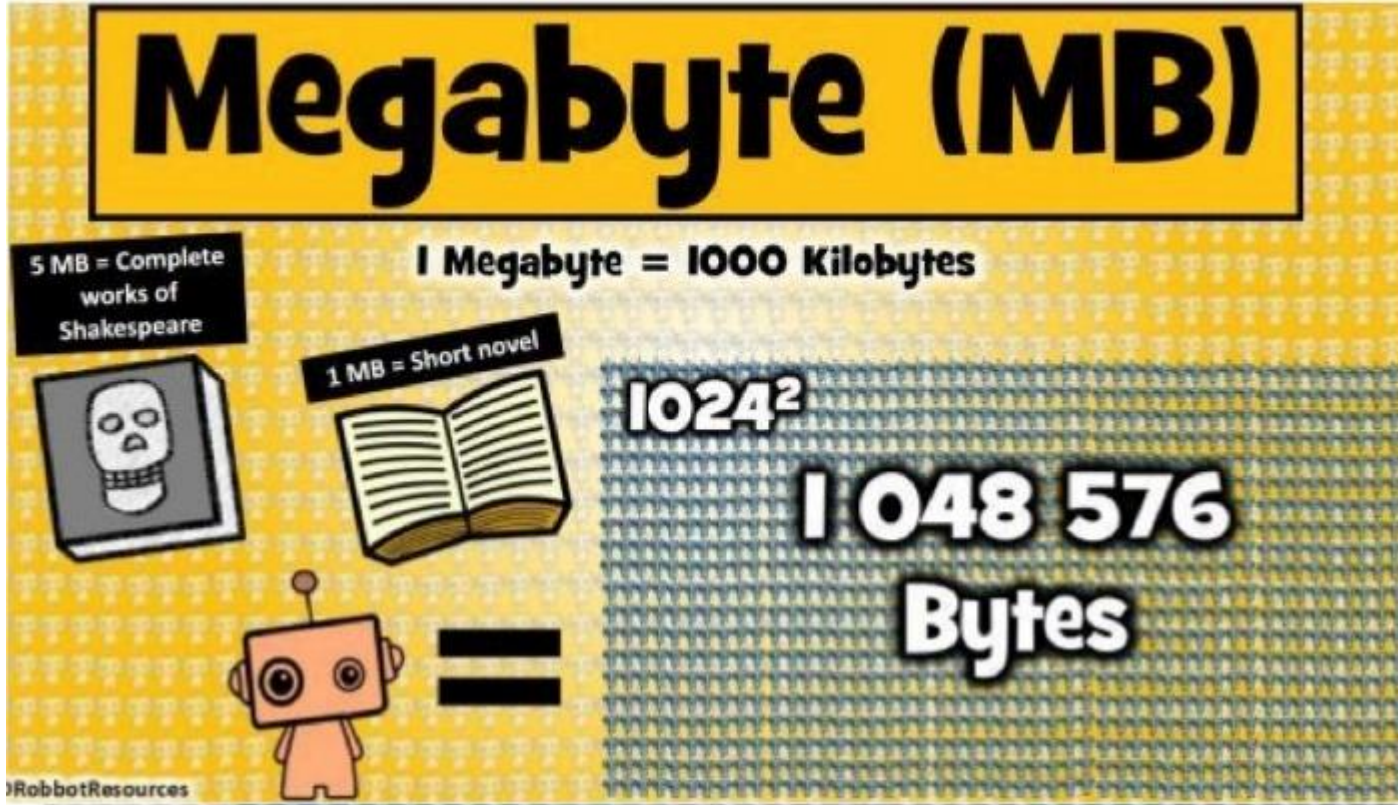
1 024 Bytes

1 KB = paragraphe

100KB = photo basse résolution

128 KB = mémoire d'une calculatrice standard

# Pourquoi utiliser le Big Data ?



1MB = roman court

5 MB = œuvres complètes de Shakespeare



# Pourquoi utiliser le Big Data ?



1GB = 7 minutes d'une vidéo HD

4.7 GB = Taille standard d'un DVD-R

16 GB = mémoire d'un smartphone moyen

# Pourquoi utiliser le Big Data ?



1TB = Disque dur d'un ordinateur portable moderne

10TB = assez pour stocker tout ce que vous regardez pendant un an

1TB = 50 000 arbres en papiers imprimés



# Pourquoi utiliser le Big Data ?



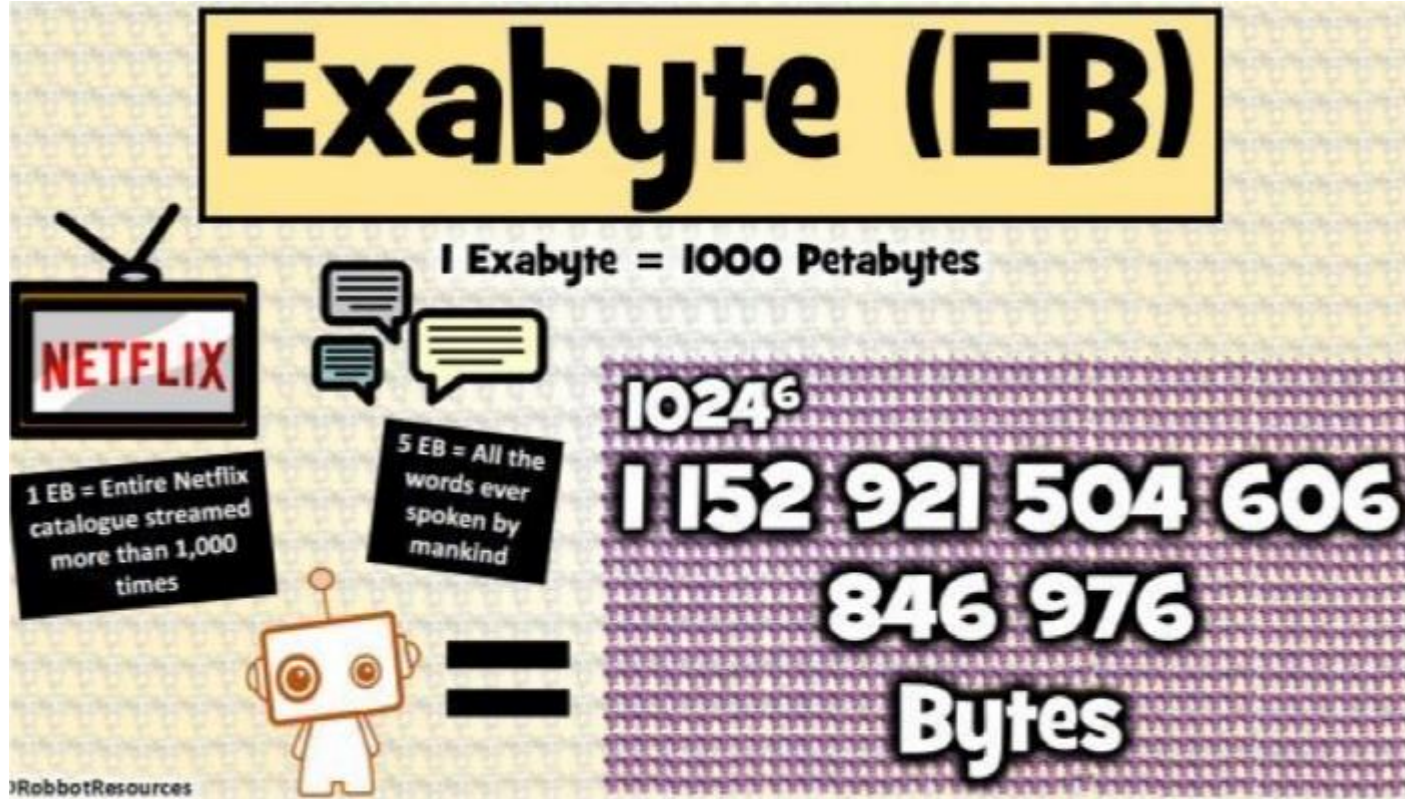
1PB = 80 millions de classeurs remplis de texte.

Le super supercalculateur Titan a coûté 93 millions de dollars et a une capacité de mémoire de 40 PB.

1.5 PB = Toutes les photos sur Facebook

20 PB = La Quantité de données traitées par Google chaque jour.

# Pourquoi utiliser le Big Data ?



1EB = l'ensemble du catalogue Netflix diffusé plus de 1000 fois.

5EB = Tous les mots jamais prononcés par l'humanité.



# Pourquoi utiliser le Big Data ?



1ZB = Chaque jour, en moyenne, environ 500 000 000 tweets sont tweetés sur Twitter. Si ces données sont stockées, il faudrait environ 100 ans pour égaler un zettaoctet de données.

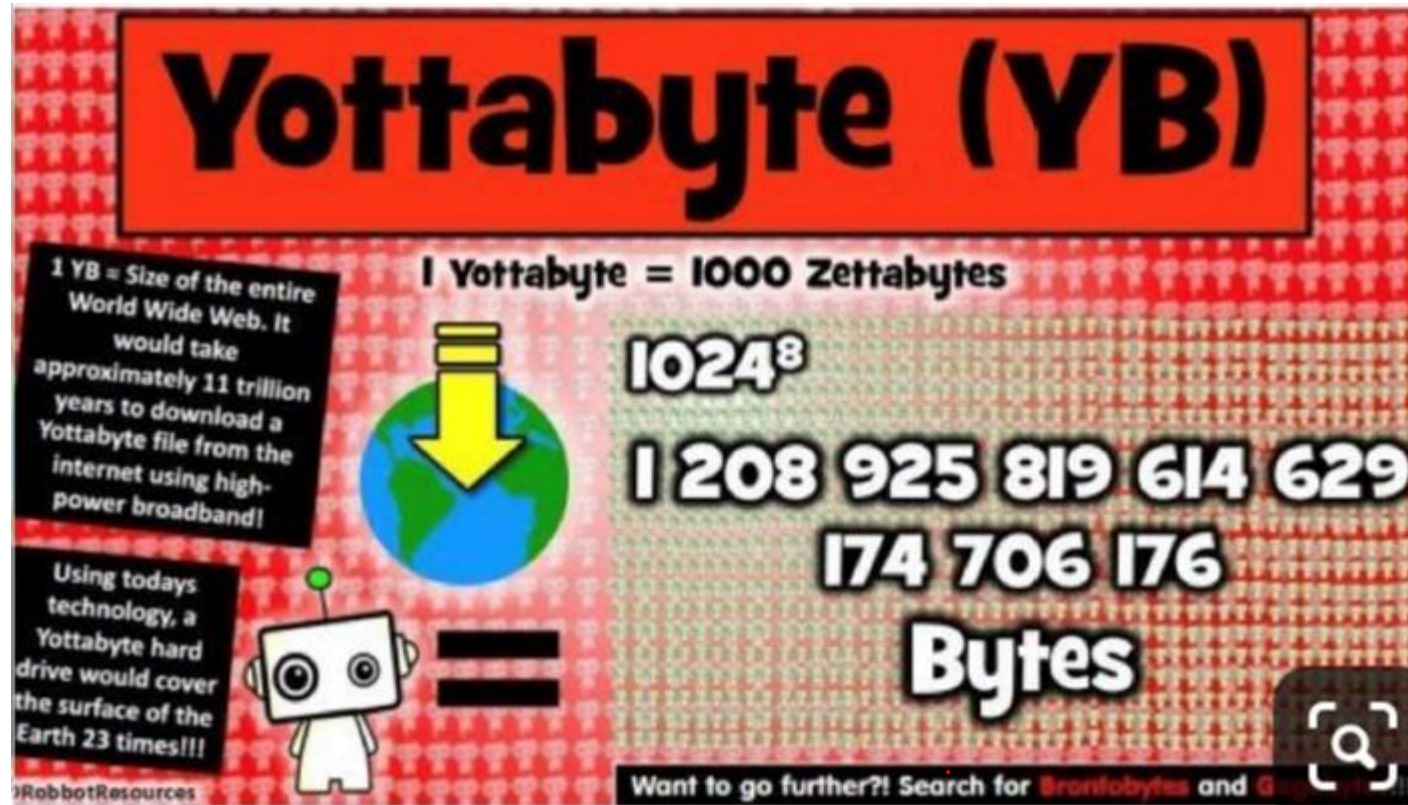
1ZB = 250 billion DVD

En utilisant la technologie d'aujourd'hui, un disque dur Zettabyte aurait la taille de l'Antarctique.



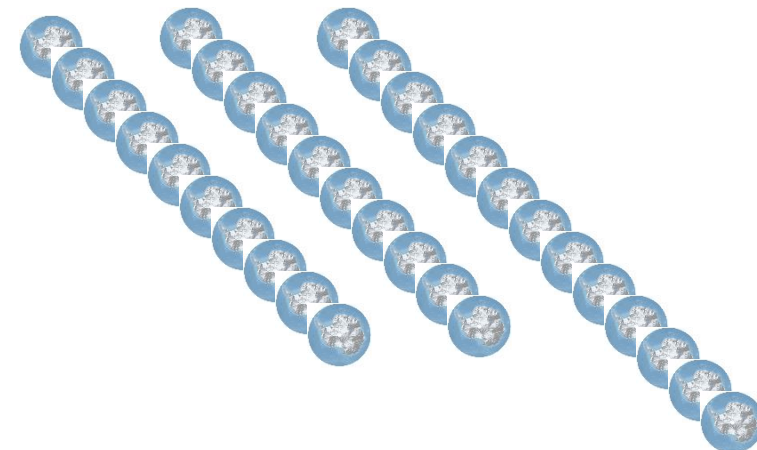


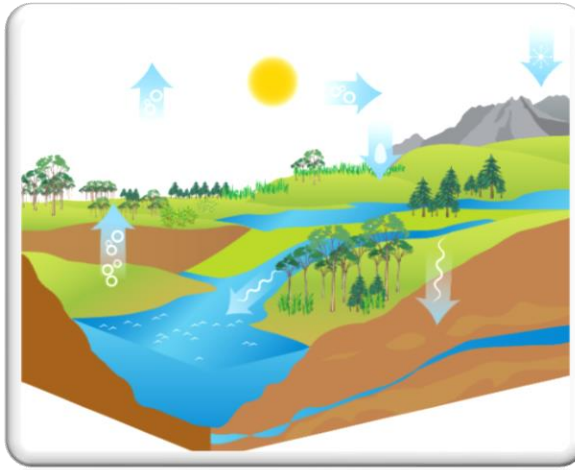
# Pourquoi utiliser le Big Data ?



1YB = Taille du World Wide Web. il faudrait environ 11000 milliards d'années pour télécharger un fichier Yottabyte à partir d'Internet à l'aide d'une bande à haute puissance.

En utilisant la technologie d'aujourd'hui, un disque dur de Yottabytes couvrirait la surface de la terre 23 fois !!!





# L'écosystème du Big Data



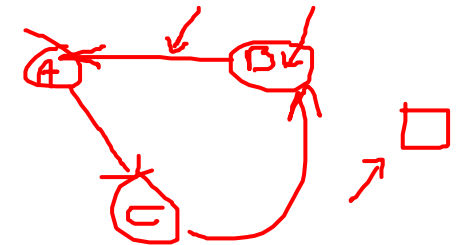
## Quantcast File System (QFS)



# Bases NoSQL



elasticsearch





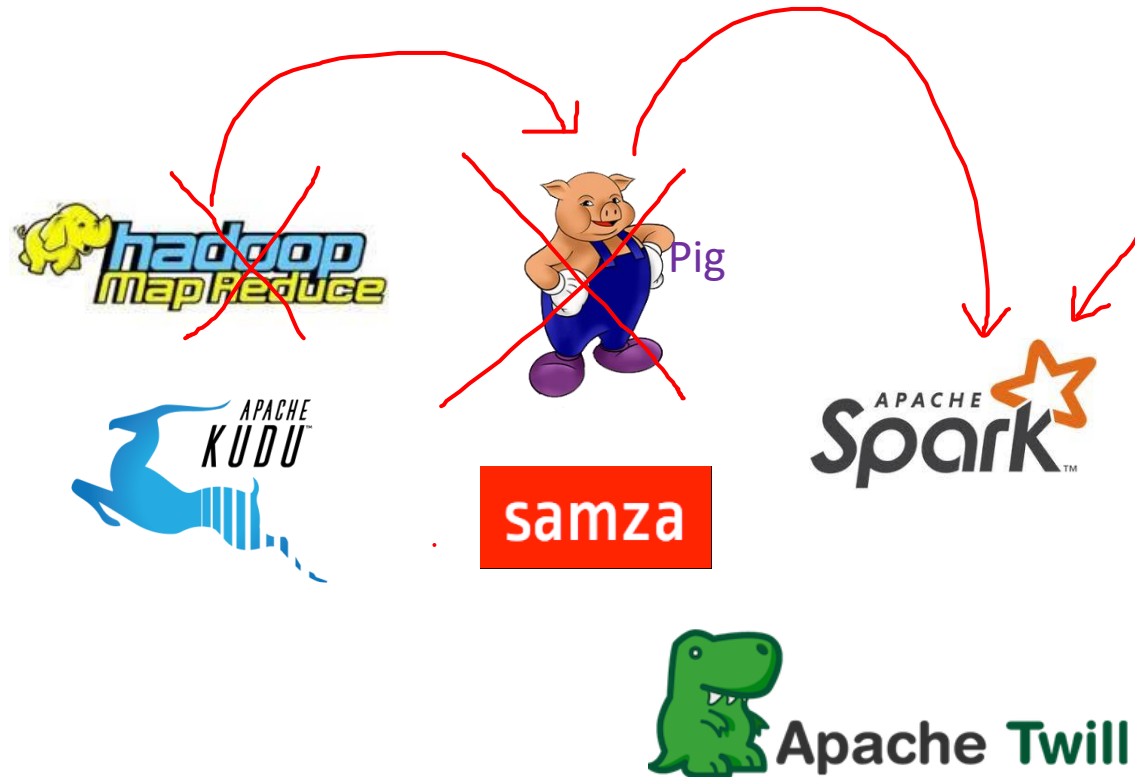


# Intégration de données



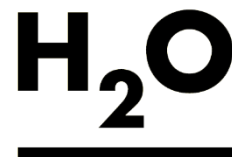


# Programmation distribuée





# Machine Learning











# L'écosystème du Big Data

## Visualisation



## Analyses des données



## Moteurs Recherche



## Récupération et ETL



## Temps Réel



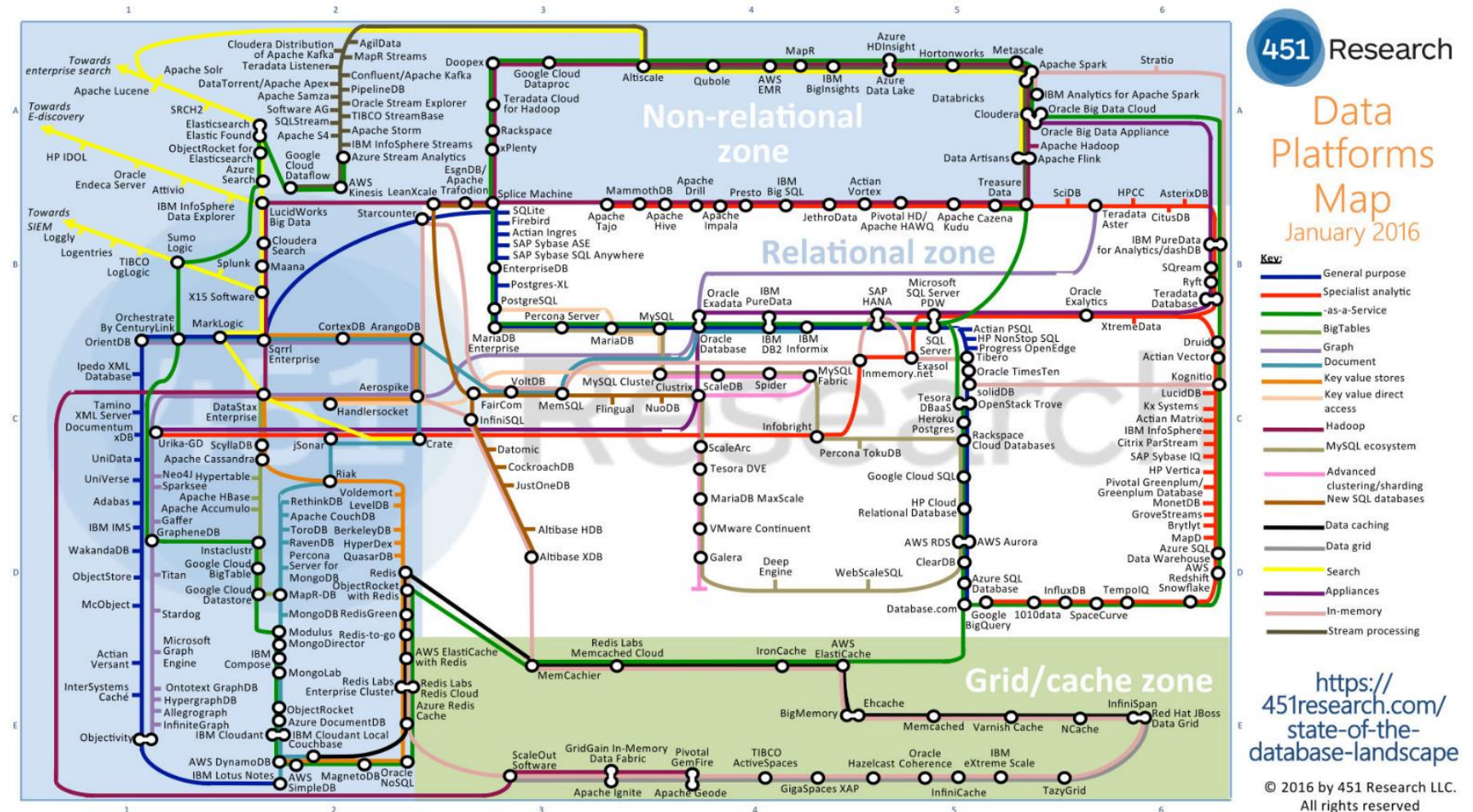
## Stockage



IllustraData.com



# L'écosystème du Big Data





# Sommaire

1. Généralités
2. **Les opportunités qu'offre le Big Data pour la DSI**
3. Open Data
4. Cas d'usage du Big Data
5. Le stockage dans le Big Data
6. Les technologies du Big Data
7. Le traitement des données en Big data
8. Compétences autour du Big Data
9. Les étapes d'un projet Big Data







# Pour le Big Data l'entreprise est la partie prenante

Plus les informations sont de qualité et plus les décisions prises sur la base de ces dernières seront pertinentes, créant ainsi de la valeur pour l'entreprise.



➤ data mining/analytics

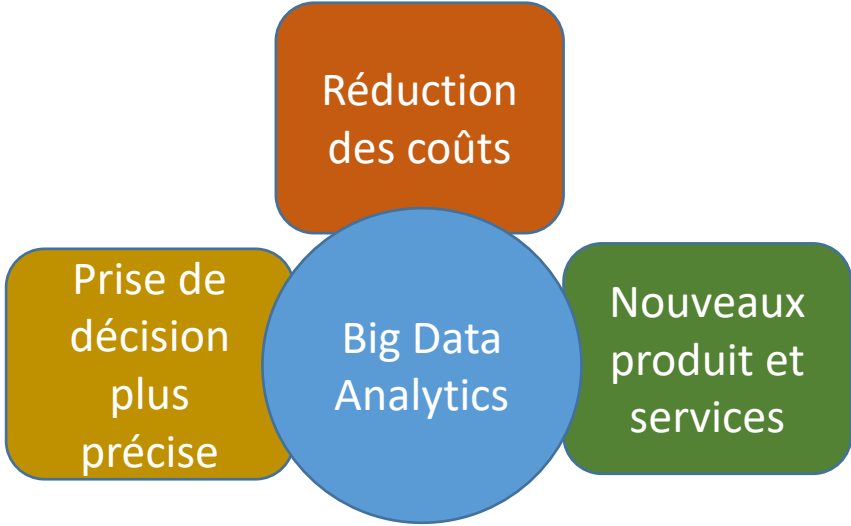
 **Big Data Analytics**





- ✓ des schémas cachés,
- ✓ des corrélations inconnues,
- ✓ des tendances du marché,
- ✓ des préférences client

**Résultat** : aider les entreprises à prendre des décisions plus efficace et avec de la **valeur**.





# La transformation des processus métier

S'assurer de la qualité des données est un moyen de limiter les risques

Il est recommandé que les données soient générées et traitées par les processus métier

Il est aussi recommandé d'identifier dans un premier temps les parties prenantes et les enjeux

L'urbanisation de la donnée (suivre l'évolution des données)

Pour urbaniser il faut gouverner ses données



# Les impacts du Big Data

# Le Big Data peut impacter l'entreprise de différentes manières

**Gouvernance** : quelles données inclure, et comment définir et assurer la gouvernance du Big Data ?

**Planification** : comment collecter les données? Comment l'analyser ? Pour quelles finalités ?  
Comment organiser les résultats ?

**Utilisation** : quelle infrastructure mettre en œuvre ? Une solution cloud serait la plus adaptée ?

**Qualité des données** : comment assurer la qualité des données traitée dans le cadre d'un certain nombre de domaines tels que la normalisation, l'harmonisation et la rationalisation ?

## Vie privée : comment protéger la vie privée ?



# RGPD

Le **RGPD** est le **R**èglement **G**énéral sur la **P**rotection des **D**onnées, une nouvelle réglementation européenne entrée en vigueur le **25 mai 2018**.

C'est le règlement no 2016/679, qui constitue le texte de référence en matière de protection des données à caractère personnel. Il renforce et unifie la protection des données pour les individus au sein de l'Union européenne.

Cette nouvelle loi a différents objectifs :

- Renforcer les droits des personnes ;
- Responsabiliser les acteurs traitant des données;
- Crédibiliser la régulation.



# Données personnelles

Les données personnelles se répartissent dans trois catégories

- **Données fournies volontairement** : il s'agit des données créées et partagées de manière explicite par des individus (par exemple, les profils sur les réseaux sociaux).
- **Données observées** : il s'agit des données collectées lors du suivi d'actions effectuées par des individus (par exemple, les données de géolocalisation issues de l'utilisation des téléphones mobiles).
- **Données déduites** : il s'agit des données sur les individus résultant de l'analyse des données des deux premières catégories.





## Ce qu'il faut retenir :

- 72









# Sommaire

1. Généralités
2. Les opportunités qu'offre le Big Data pour la DSI
- 3. Open Data**
4. Cas d'usage du Big Data
5. Le stockage dans le Big Data
6. Les technologies du Big Data
7. Le traitement des données en Big data
8. Compétences autour du Big Data
9. Les étapes d'un projet Big Data



# Donnée (Data)

1640s : première utilisation anglaise du mot "data"

1946 : "données" signifie "informations informatiques transmissibles et stockables".

1954 : première utilisation anglaise du terme "data processing".

Une donnée est une description élémentaire d'une réalité.

- par exemple une observation ou une mesure







# Codage de l'information

## Exemples :

## Code à barres



## QR Code



# Radio Frequency IDentification





# Atelier

Chaque participant doit trouver un sujet Big Data à réaliser par un prestataire afin de créer une valeur.

- Définir un sujet Big Data
- Les données récoltées (!!! RGPD !!!)
- La valeur de rentabilité

## Données externes non-structurées:

Réseaux, UGC, etc.



## Données externes structurées:

Panel, Sondage, Partenaires, etc.



person	year	income	age	sex
1	2001	1300	27	1
1	2002	1600	28	1
1	2003	2000	29	1
2	2001	2000	38	2
2	2002	2300	39	2
2	2003	2400	40	2

## Données internalisées:

Datawarehouse, ERP, CRM, etc.



Données utilisables par  
l'entreprise pour mettre en place  
des solutions

