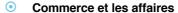


Aucun domaine n'échappe à l'avalanche des données



- SI d'entreprise, Banques, transactions commerciales, systèmes de réservation, ...
- Gouvernements et organisations
 - O Lois, réglementations, standards, infrastructures,
- Loisirs
 - O Musique, vidéo, jeux, réseaux sociaux...
- Sciences fondamentales
 - O Astronomie, physique et énergie, génome, ...
- Santé
 - O Dossier médical, sécurité sociale,...
- Environnement
 - O Climat, dév durable, pollution, alimentation,...
- Humanités et Sciences Sociales
 - Numérisation du savoir (littérature, histoire, art, architecture), données archéologiques...

Mokrane Bouzeghoub

5

Une grande variété de données

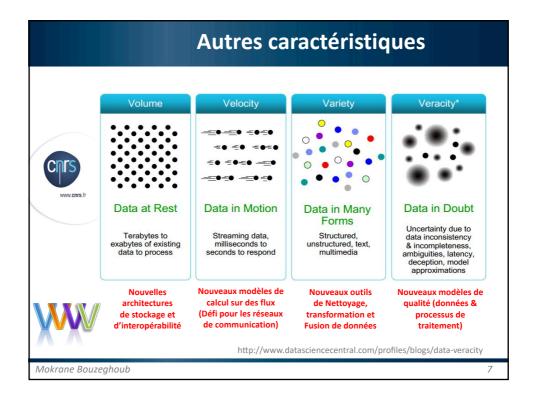
- Données de références
 - O Cadastre, SIG, catalogues de produits, ...
- Données d'observations
 - O Satellites, capteurs, expériences scientifiques, ...
- Données transactionnelles
 - O Transactions commerciales, requêtes BD/Web, ...
- Données sociales
 - Web, Facebook, Twitter, Crowdsourcing, ...
- Données du patrimoine
 - Culture, architecture, publications, ...
- ...

Représentées sous forme

de tables, de graphes, de textes, d'images ou de flux vidéo

Mokrane Bouzeghoub





Big Data – Définitions différentiées



• Chercheur informaticien: c'est la frontière au delà de laquelle les outils dont je dispose ne permettent pas le traitement effectif des données, en raison de leur volume ou en raison de méthodes inadaptées

Non-informaticien : C'est la mise en œuvre de nouvelles méthodes de traitement et d'analyse de données susceptibles de créer une rupture dans la façon de pratiquer le métier d'analyste ou de chercheur

Mokrane Bouzeghoub

Caractéristiques du domaine

- Un domaine très vaste,
 - en interaction permanente avec les autres disciplines scientifiques



- Un domaine qui se repositionne périodiquement
 - En revisitant ses solutions à la lumière de nouvelles technos et de nouvelles idées
 - En intégrant de nouveaux besoins et de nouveaux problèmes
- Une recherche dominée (ou presque) par des labos industriels :
 - Google, Facebook, Yahoo!, Amazone, IBM, Oracle, Microsoft ...

Mokrane Bouzeghoub

9

Les grands challenges du Big Data

- Stockage et préservation des données
 - O Performance des accès, disponibilité des données
 - O Protection des données, qualité des données
 - Indexation sémantique (ontologies), indexation participative (folksonomies)



- Analyse statistique et sémantique, raisonnement
 - Analyse en temps réel de flux continus de données émanant de différentes sources
 - Requêtes multidimensionnelles sur des grands ensembles de données
 - Extraction et interprétation de connaissances, apprentissage profond
- Impact social et économique
 - O Protection de la vie privée, Droit à l'oubli
 - Droits de propriétés, droits d'exploitation
 - O Economie d'énergie, écologie
 - coût du stockage, coût de transfert



Mokrane Bouzeghoub

Le Big Data n'est pas gratuit

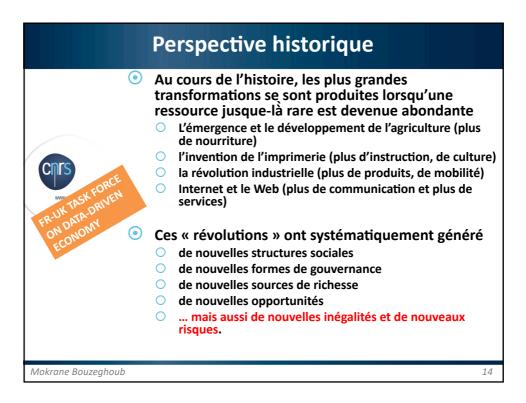
- Consommation électrique des data centers
 - Entre 600 TWh et 1 500 TWh, soit environ 10 % de l'électricité produite dans le monde, équivalent de la production de 180 réacteurs nucléaires
- Coût d'extraction des matières premières (terres rares)
 - on non durable et souvent non éthique
- Coût de production et de renouvellement des équipements
 - Accélération de l'innovation
- Coût de recyclage des déchets
 - O Un cycle partiellement maîtrisé

Mokrane Bouzeghoub

11

La recherche en big data est un continuum de la recherche sur les données **Plusieurs** Stockage, indexation, distribution, data décennies de recherche et Requêtes continues, requêtes d'innovation approximatives Requêtes avec préférences, requêtes skyline **VLDB** Analyse de flux, agrégation en ligne, algèbre OLAP **EDBT** Very Big Data Intégration de données, systèmes de médiation. ETL Fouille de données, découverte de motifs Exécution distribuée et parallèle de requêtes Data Deluge Big Data Mokrane Bouzeghoub





Spécificité de la ressource « donnée »

On doit investir
dans les données
de la même
manière que nous
investissons dans
les routes, les
chemins de fer et
d'autres
infrastructures
publiques.
[FR-UK Report on
pata-Driven Economy]

- l'importance stratégique des données n'est pas encore reconnue par tous
 - O Peuvent faire l'objet d'un usage partagé
 - Ne sont pas détruites après utilisation
 - O La production de données est quasi infinie

les données n'ont pas toujours de valeur intrinsèque

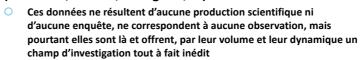
- Leur valeur résulte des services, produits et connaissances qui génèrent de nouveaux modèles économiques et de nouvelles activités
- leur utilisation répétée peut augmenter leur valeur et conduire à de nouvelles formes de ressources (métadonnées, connaissances)

Mokrane Bouzeghoub

15

De nouveaux types de traitements sur les données







- Analyse d'opinions, identification de leaders, détection de signaux faibles, recommandations commerciales, ...
- À la portée du plus grand nombre
- ... Et un usage réflexif par les algorithmes
 - amélioration des algorithmes d'apprentissage et de recherche d'information, en faisant un retour d'expérience à ces algorithmes.

Avec de nouveaux outils

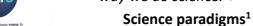
Fouille de données, apprentissage, visualisation, sciences participatives (annotations,...)

Mokrane Bouzeghoub

Précurseur de la science des données

N'oublions pas que cette problématique a été mise en lumière par Jim Gray, en janvier 2007, disparu en mer deux semaines plus tard.

« The availability of very large amounts of data and the ability to efficiently process them is changing the way we do science. »



- I. Empirical description of natural phenomena
- 2. Theoretical science: models and generalization
- 3. Computational science: simulation of complex phenomena to validate theories
- 4. Data Intensive science : collecting and analyzing large data volume

 1 Jim Gray, eScience Talk at NRC-CSTB meeting Mountain View CA, 11

Turing Award



Mokrane Bouzeghoub

17

Science des données et l'infiniment petit Chromosomes Génome Cellule Groupe de cellules Chromosomes Génome Cellule ADN Cellule vivante Protéines Pour observer et comprendre le Vivant Génétique, Phylogénie, Biologie évolutive... Pour comprendre la physique des particules LHC, énergie Mokrane Bouzeghoub



Science des données dans le quotidien à l'ère du Net Internet a initié une nouvelle ère commerciale En utilisant nos données personnelles et nos interactions (Les likes, les tag, les préférences, Les données personnelles) En exploitant ces données pour de nouvelles opportunités commerciales Les réseaux sociaux deviennent un observatoire de la Société Qui parle de quoi ? Quels sont les thèmes de discussion majeurs? Quels sont les thèmes émergents? Qui sont les leaders d'opinion? Peut-on détecter des signaux faibles ? · Suicide, attentat, pédophilie, Mokrane Bouzeghoub

Science des données dans l'économie et le commerce



- En 2014, un rapport du McKinsey Global Institute a conclu, que le Big Data pourrait rapporter 3 000 Mds \$ chaque année dans seulement 7 secteurs d'activité
- En 2015, la Commission européenne a estimé que la taille cumulée du marché direct des données publiques ouvertes devrait s'élever à 325 Mds € en Europe pour la période 2016-2020
- En 2015, 61 % des sociétés françaises estimaient que le Big Data était devenu l'un des principaux moteurs de croissance, aussi important pour elles que leurs produits et services existants
- en 2018, le marché français du Big Data devrait atteindre 652 M€, soit une hausse de 129 % par rapport à 2014, malgré une conjoncture défavorable

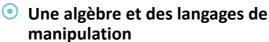
Mokrane Bouzeghoub

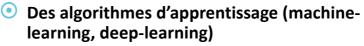
21

Quelles techniques pour faire parler les données ?







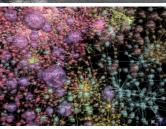


- Des algorithmes de fouille de texte
- Des algorithmes d'analyse d'images
- Des algorithmes de traitement du signal (flux audio ou vidéo)
- ...

Mokrane Bouzeghoub

Mais très souvent ... seulement une exploration visuelle





visualcomplexity.com/vo

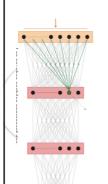
 Pour observer des phénomènes non perceptibles à l'œil nu ou à l'aide d'instruments

- durant la simulation par exemple
- O Interaction entre faisceaux de particules
- Pour analyser à un niveau macro des interactions entre objets
 - Phylogénie
 - Mouvement de foule, inondations
 - Réseaux sociaux
- Attention
 - Si la visualisation peut aider à la compréhension d'un phénomène, elle peut introduire un biais et en altérer l'interprétation?
 - → Recherches fondamentales sur l'interaction H-M

Mokrane Bouzeghoub

23

L'apprentissage automatique : l'arbre qui cache la forêt



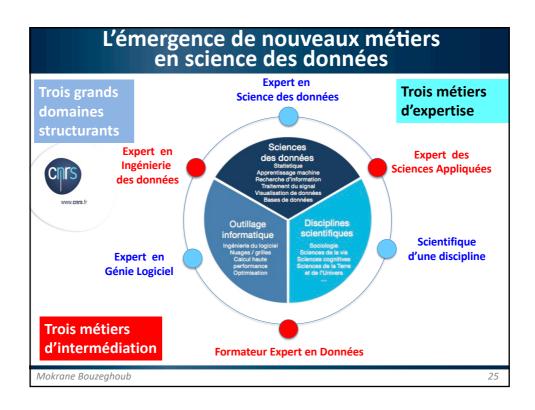
Un outil d'extraction de connaissances et de reconnaissances de formes très performant

- Nécessitant des millions d'exemples annotés qui vont servir pour l'entrainement
- Basé sur des algorithmes très complexes et très gourmands en ressources de calcul
- Mais n'intervient qu'en bout de la chaîne de valorisation des big data
 - Acquisition des données
 - Qualification, annotation, nettoyage des données
 - Intégration des données
 - Stockage, indexation, tri, accès optimisé
 - 0 .

Pas pertinent pour tous les besoins

- La corrélation de certaines données peut se faire par de "simples jointures" ou quelques opérateurs analytiques ("skyline")
- Un bon langage de requêtes suffit pour beaucoup d'applications

Mokrane Bouzeghoub





Objectifs oissante du n

- La complexité croissante du monde ne peut plus se satisfaire d'une approche scientifique unique
 - éclatement de la connaissance
 - O pluralité des savoirs sur des mêmes faits sociaux et sur leur articulation



- Favoriser l'émergence d'une communauté scientifique interdisciplinaire autour de la science des données, et produire des solutions originales sur le périmètre des données scientifiques.
- Produire des concepts et des solutions qui n'auraient pu être obtenus sans coopération entre les différentes disciplines

Mokrane Bouzeghoub

27

Indicateurs de Suivi

- Pérennité de la coopération
- Publications communes



- Co-encadrement de thèses
- Plateformes de test et d'expérimentation
- Colloques/journées de dissémination
- Montage et soumission de nouveaux projets
- Synergie entre projets



Mokrane Bouzeghoub

Quelques chiffres

- Défi lancé en 2012
 - O 50 projets retenus sur 120 soumis
- Budget



- Près de 3,5 M€
- O Montant alloué/projet/an 20 à 120 K€
- Degré d'implication des unités
 - Plus de 100 labos impliqués, couvrant les 10 instituts du CNRS
 - O Plusieurs EA universitaires associées
 - O Plus de 400 chercheurs impliqués
- Organismes impliqués hors CNRS



INRIA, INRA, IRSTEA, INSERM, CEA, ONERA, Universités, Ecoles

Mokrane Bouzeghoub

29

Thématiques couvertes par les AAP

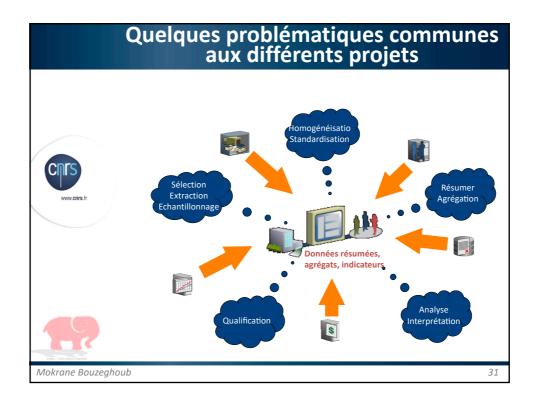
- Collecte, stockage et indexation de données massives
- Hétérogénéité, interopérabilité, intégration, partage des données

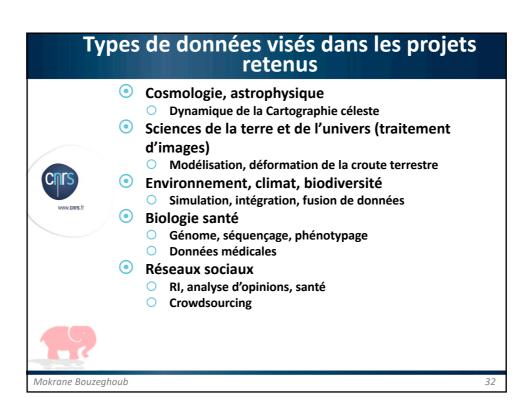


- Calcul intensif sur des grands volumes de données, parallélisme dirigé par les données, optimisation
- Extraction de connaissances, datamining et apprentissage
- Préservation/archivage des données (animation)
- RI agrégative/résumé, sémantique et raisonnement
- Visualisation de grandes masses de données
- Qualité, protection et sécurité des données
- Problèmes de propriété, de droit d'usage, droit à l'oubli
 - Consommation d'énergie, environnement, recyclage...



Mokrane Bouzeghoub





Bilan à 6 ans

- Une communauté interdisciplinaire en formation sur les thèmes des Big Data et Science des Données
 - Organisée dans un GdR d'animation (MaDICS, plus de 600 adhérents en déc. 2015)



- Une production scientifique interdisciplinaire
 - Plusieurs publications signées par des chercheurs de différentes disciplines
 - Plusieurs co-encadrements de thèses
- Nouvelles initiatives de recherche issues du défi
 - O Défi Imag'In sur le traitement d'images
 - O Le PEPS FaSciDo sur les fondements des sciences de données
 - Le PEPS Astro-Informatique sur les sciences de données en astrophysique
- Levier pour initier d'autre projets ou réseaux de compétences
 - O Plus de 10 projets ANR



2 actions COST (Europe)

Mokrane Bouzeghoub

2:

Le GdR MaDICS

Actions actuelles

- Apprentissage, optimisation Large-échelle et calculs distribués (ATLAS)
- 2. Qualité des données scientifiques (ARQUADS)
- 3. Environmental Acoustic Data Mining (EADM)
- 4. Graph Data Mining (GRAMINEES)
- 5. Imagerie Hyperspectrale (Imhyp)
- 6. Masses de données en astronomie et astrophysique (MAESTRO)
- 7. Préservation des données scientifiques (PREDON)
- 8. Reproductibilité des expériences d'analyse de données scientifiques (ReProVirtuFlow)
- 9. Raisonner sur les données (RoD)



http://www.madics.fr

Mokrane Bouzeghoub

