

Le déluge de données comment en tirer parti

Jeudi 9 juin 2011

Coordination scientifique :

- *David Menga (EDF R&D)*
- *Jean-Michel Batto (INRA)*
- *Pierre Léonard (INRA)*

Amphithéâtre Becquerel, École Polytechnique, Palaiseau

<http://www.association-aristote.fr>

info@association-aristote.fr

Edition du 18 prairial an CCXIX (*vulg.* 7 juin 2011) ©2011 Aristote

Table des matières

1	Programme de la journée	5
1.1	Introduction	5
1.2	Programme	6
2	Présentations	7
	Pierre Léonard & Jean Michel Batto (INRA), David Menga (EDF R&D)	7
2.1	Xavier DALLOZ (Consultant)	9
2.2	Vincent HEUSCHLING (D-FI)	16
2.3	Alexandru Costan (INRIA)	21
2.4	Joerg Bienert (ParStream)	30
2.5	Ronan Keryell (HPC-Project)	37
2.6	Denis Caromel (ActiveEon-INRIA)	56
2.7	Nicolas Pons (INRA)	61
2.8	Patrick Furhmann (DESY-Hamburg)	69
2.9	Marie-Luce Picard (EDF R&D et ENST-Bilab)	76
2.10	David KONOPNICKI (IBM Haïfa-Research)	82

Chapitre 1

Programme de la journée

1.1 Introduction

De nos jours, les données sont partout et constituent la matière première de notre monde numérique. Elles redéfinissent la façon dont on crée de la connaissance scientifique et offrent aux entreprises de nouveaux leviers de croissance et plus de performances opérationnelles aux gouvernements.

Décoder le génome signifie analyser 3 milliards de paires de base. Walmart traite chaque heure 1 million de transactions clients avec une volumétrie totale de 2,5 Petaoctets. Nous vivons « une révolution industrielle » des données et nos outils actuels de stockage et de traitement des données sont inadaptés pour traiter de tels volumes en un temps acceptable.

Ce séminaire a pour objet de donner les clefs de compréhension de cet univers, de mettre en perspective les défis scientifiques et opérationnels et d'offrir des éléments de réponse à travers des travaux scientifiques de pointe et des solutions marché innovantes.

Le matin sera consacré à la définition de la problématique et donnera des réponses. L'après midi, nous aborderons les expériences des acteurs confrontés à ce défi, comme les électriciens, les biologistes et les astrophysiciens.



Photos ©2011 Philippe Laviolle

1.2 Programme

9h00-9h30	<i>Accueil café</i>	
9h30-9h45	Pierre Léonard (INRA) Jean Michel Batto (INRA) David Menga (EDF R&D)	Présentation du séminaire
9h45-10h30	Xavier Dalloz (Consultant)	<i>Big data</i> , le fuel de l'économie du XXI ^e siècle
10h30-11h00	<i>Pause café</i>	
11h00-11h45	Vincent Heuschling (D-FI)	Cartographie des solutions <i>big data</i> du marché
11h45-12h30	Alexandru Costan (INRIA)	Analyse des systèmes de stockage à grande échelle pour les applications de traitement intensif des données
12h30-14h00	<i>Repas (salle Detoef)</i>	
14h00-14h35	Joerg Bienert (ParStream)	An innovation solution to manage heterogeneous big data
14h35-15h10	Ronan Keryell (HPC-Project)	Environnement de programmation pour traitements massifs sur architectures modernes
15h10-15h20	Denis Caromel (ActiveEon-INRIA)	Solutions ProActive pour <i>Workflows, Map/Reduce, Matlab /Scilab, CPU/GPU</i>
15h20-15h45	Nicolas Pons (INRA)	La métagénomique, un défi supplémentaire pour la loi de Moore
15h45-16h00	<i>Pause</i>	
16h00-16h35	Patrick Furhmann (DESY-Hamburg)	dCache : scaling out affordable storage.
16h35-17h10	Marie-Luce Picard (EDF R&D et ENST-Bilab)	Données massives pour les <i>smart-grids</i>
17h10-17h40	David Konopnicki (IBM-Haïfa Research Labs)	Massive Scale Analytics for a Smarter Planet
17h45-18h00	Conclusions, questions-réponses avec les intervenants	

Chapitre 2

Présentations

Pierre Léonard & Jean Michel Batto (INRA), David Menga (EDF R&D)

Ouverture du séminaire



Big Data Le Déluge de données, Comment en tirer parti

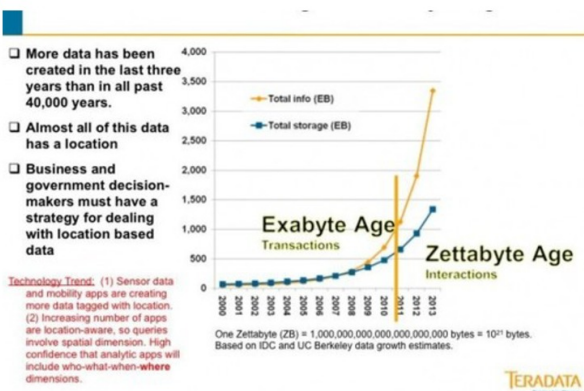
Jean-Michel Batto, INRA
Pierre Léonard, INRA
David Menga, EDF R&D

Séminaire Aristote, 9 juin 2011

1



2



Séminaire Aristote, 9 juin 2011

3

Les objets produisent des data en temps réel

- Un seul moteur de Boing produit 10 To d'informations toutes les 30 secondes
- Un voyage de 6 heures New-York/Los Angeles en 737 produit 240 To
- Chaque jour, il y a 28537 vols commerciaux aux Etats Unis



1 Pétaoctet / Jour

Le StoriQ ST32R-1PT, serveur SAN d'Intelligence, 1 Pétaoctet, coûte 300,000 euros, soit moins de 30 centimes/Go

Séminaire Aristote, 9 juin 2011

4

De nouveaux concepts, BI et BA

- La Business Intelligence ou Informatique décisionnelle désigne les outils et méthodes pour collecter des données et mesurer les performances d'un système à travers des indicateurs bien choisis. Le BI répond aux questions :
 - Que s'est-il passé ?
 - Où sont les problèmes ?
 - Que doit-on faire ?
- Le Business Analytics (BA) fait référence aux technologies et méthodes pour analyser des données afin de comprendre ce qui se passe et de prédire les évolutions futures. Le BA répond aux questions
 - Pourquoi est-ce arrivé ?
 - Que se passera-t-il plus tard si la tendance perdure ?
 - Comment optimiser le système ?

Séminaire Aristote, 9 juin 2011

5

14h00-14h35	Joerg Bienert (Parstream)	An innovative solution to manage heterogeneous big data
14h35-15h10	Emmanuel Chua (HPC-Project)	Que l'environnement de programmation pour traiter les big data
15h10 -15h 20	Denis Caronnel (ActinEon-INRIA)	Solutions ProActive pour Workflows MapReduce, Metadab/SciLab, CPU/GPU
15h 20-15h45	Nikolas Pons (INRA)	La mégéconomiq., un défi supplémentaire pour la loi de Moore
15h 45- 16h 00	Pause	
16h 00-16h35	Patrick Fuhmann (DESY)	dCache: scaling out affordable storage
16h35-17h10	Marie-Luce Picard (EDF R&D et ENST-BIB)	Données massives pour les Smart-grids
17h 10- 17h45	David Konopnicki (IBM Haifa-Research Lab)	Massive Scale Analytics for a Smarter Planet
17h45-18h00	Conclusion, questions réponses avec les intervenants	

Séminaire Aristote, 9 juin 2011

7

Programme

09h00-09h30	Accueil café	
9h30-9h45	Pierre Léonard, Jean Michel Batto (INRA), David Menga (EDF R&D)	Présentation du Séminaire
9h45-10h30	Xavier Daltoz (Consulting)	<i>Big data</i> , le fuel de l'Economie du XXI ^{ème} siècle
10h30-11h00	Pause café	
11h00-11h45	Vincent Heuschling (D-FI)	Cartographie des solutions <i>big data</i> du marché
11h45-12h30	Alexandru Costan (doctorant INRIA)	Quel système de gestion de données pour stocker les données massives ?
12h30-14h00	Déjeuner (salle Deteouf)	

Le Domesday Book (1086)

- Le **Domesday Book** (ou simplement **Domesday**), en français *Livre du Jugement Dernier*¹, est l'enregistrement du grand inventaire de l'**Angleterre** terminé en **1086**, réalisé pour **Guillaume le Conquérant**, l'équivalent de nos jours d'un **recensement national**.



Séminaire Aristote, 9 juin 2011

8

2.1 Xavier DALLOZ (Consultant)

***Big data*, le fuel de l'économie du XXI^e siècle**

Régulièrement une innovation majeure change tout avec à chaque fois de nouvelles technologies, de nouveaux métiers et de nouveaux enjeux. Après l'ère des *mainframes*, l'ère des *mini computers*, l'ère des PC, l'ère du *software* pour améliorer la productivité personnelle, l'ère de l'Internet, voici celle du « *Big Data* ».

Le *Big Data* va permettre notamment de repenser les modèles économiques en misant sur de nouvelles créations de valeur avec notamment l'intelligence collaborative et la *shazamisation* de notre environnement de telle sorte qu'il y ait une meilleure efficacité de « notre » capital : santé, énergie, éducation, équipements, stocks... Avec le *Big Data*, la chasse aux gaspillages va enfin devenir une réalité. Nous n'avons encore rien vu... Tout va s'accélérer.

Que d'opportunités ! Que de leviers de croissance pour nos économies !

Les datas = le fuel du 21ème siècle

D'énormes gisements de création de valeurs



<http://www.yournetworkmarketing.com/facebook-twitter-youtube-stats-in-real-time-simulation/>

Xavier Dalloz

dalloz@dalloz.com

XAVIER DALLOZ
CONSULTING

Le Plan

- Définition des Big Data : un réducteur de la complexité + une aide à la prise de la décision et au management du risque (intuition)
- Exemples de Big Data
- Les outils du Big Data
- L'enjeu : la création de la valeur en optimisant l'utilisation du capital fixe, circulant et immatériel
- Exemples d'usage des Big Data
 - Le commerce de détail (stocks)
 - La production des biens manufacturés (co-création)
 - La géolocalisation (traces)
 - L'Internet des objets et la shazamisation (conversation)
 - Le web sémantique (meta data)
- Recommandations

dalloz@dalloz.com

XAVIER DALLOZ
CONSULTING

Définition des Big Data (Wikipedia)

Une définition qui change avec le temps...

- Ensembles de données qui deviennent tellement gros qu'ils en deviennent difficiles à travailler avec des outils classiques de gestion de base de données.
- Les perspectives du traitement des big data sont énormes, notamment pour l'analyse d'opinions ou de tendances industrielles, la génomique, l'épidémiologie ou la lutte contre la criminalité.
- La production de données par les utilisateurs, et notamment le **partage d'informations ubiquitaires** (capteurs et senseurs mobiles, caméras, microphones, appareils photos, lecteurs RFID, réseaux de capteurs sans fil, etc.) augmentent drastiquement le nombre de données pouvant être traitées

dalloz@dalloz.com

XAVIER DALLOZ
CONSULTING

Une autre définition

- Le Big Data est le moyen de valider une intuition
- En validant une intuition le Big Data est un réducteur de risques.
- Pour cela, il faut traiter un déluge de données en temps réel.
- Le grand frein du Big Data est l'incapacité à définir clairement notre intuition. Il faut être problématicien.

dalloz@dalloz.com

XAVIER DALLOZ
CONSULTING

Ce déluge de données a de nouvelles propriétés

Plus de données ont été créées ces 3 dernières années que pendant les 40.000 années précédentes

- Ce sont des données **non structurées**
- Elles sont produites en **temps réel**
- Elles arrivent mondialement en **flots continus**
- Elles sont **méta taguées** mais de façon disparate (localisation, heure, jour, etc.)
- Elles proviennent de sources très **disparates** (téléphone mobile, capteurs, téléviseurs connectés, tablettes, PC fixes, PC portables, objets, machines), de façon **désordonnée et non prédictible**.

dalloz@dalloz.com

XAVIER DALLOZ
CONSULTING

Unit	Size	What it means
Bit (b)	1 or 0	Short for "binary digit", after the binary code (1 or 0) computers use to store and process data
Byte (B)	8 bits	Enough information to create an English letter or number in computer code. It is the basic unit of computing
Kilobyte (KB)	1,000, or 2 ¹⁰ , bytes	From "thousand" in Greek. One page of typed text is 2KB
Megabyte (MB)	1,000KB; 2 ²⁰ bytes	From "large" in Greek. The complete works of Shakespeare total 5MB. A typical pop song is about 4MB
Gigabyte (GB)	1,000MB; 2 ³⁰ bytes	From "giant" in Greek. A two-hour film can be compressed into 1-2GB
Terabyte (TB)	1,000GB; 2 ⁴⁰ bytes	From "monster" in Greek. All the catalogued books in America's Library of Congress total 15TB
Petabyte (PB)	1,000TB; 2 ⁵⁰ bytes	All letters delivered by America's postal service this year will amount to around 5PB. Google processes around 1PB every hour
Exabyte (EB)	1,000PB; 2 ⁶⁰ bytes	Equivalent to 10 billion copies of <i>The Economist</i>
Zettabyte (ZB)	1,000EB; 2 ⁷⁰ bytes	The total amount of information in existence this year is forecast to be around 1.2ZB
Yottabyte (YB)	1,000ZB; 2 ⁸⁰ bytes	Currently too big to imagine

The prefixes are set by an intergovernmental group, the International Bureau of Weights and Measures. Yotta and Zetta were added in 1991; terms for larger amounts have yet to be established.
Source: *The Economist*

dalloz@dalloz.com

XAVIER DALLOZ
CONSULTING

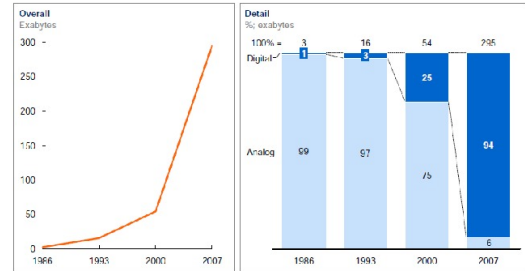
What was measured	Amount of data	Year of estimate	
MGI storage-based approach	<ul style="list-style-type: none"> New data stored in enterprise external disk storage in a year New data stored by consumers in a year 	<ul style="list-style-type: none"> 1.4 x 10¹³ bytes (includes replicas) 6.8 x 10¹³ bytes 	For 2010
IDC/EMC¹ Digital Universe	<ul style="list-style-type: none"> Annual digital data captured (includes all generated, stored or not) Includes more than 60 types of devices Did not include information consumption by users through TV, video gaming² 	- 800 x 10 ¹³ bytes	For 2009
UCSD	<ul style="list-style-type: none"> Includes both digital and analog data for TV, radio, phone, print, computer, comp games, movies, recorded music, etc. Measured data from consumption perspective² 	3.6 x 10 ¹³ bytes (total consumption US only)	For 2008
Hilbert, López	<ul style="list-style-type: none"> Capacities for specific technologies Server and mainframe hard disks Other hard disks Digital tape FC hard disk Total digital storage capacity 	<ul style="list-style-type: none"> 24.5 x 10¹³ bytes 6.49 x 10¹³ bytes 32.5 x 10¹³ bytes 123 x 10¹³ bytes 276 x 10¹³ bytes 	For 2007

1 Includes chip cards, floppy disks, camera, video games, mobiles, memory cards, media players, CDs, DVDs, Blue Ray disks, FC and server hard disks.
 2 Consumption is defined as the data each user used by the user.
 SOURCE: IDC write papers on Digital Universe, sponsored by EMC; Bolin and Shea, *How Much Information? 2005: Report on American Consumers*, January 2010; Hilbert and López, *The world's technological capacity to store, communicate, and compute information*, Science, February 2011; McKinsey Global Institute analysis

dalloz@dalloz.com

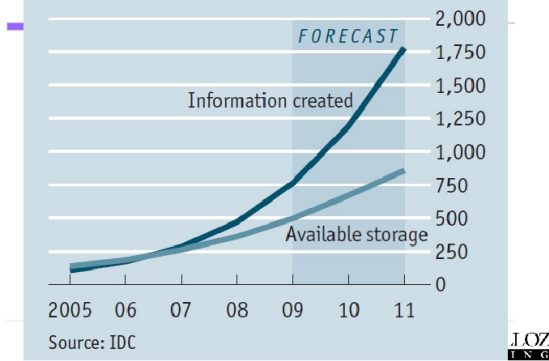
Data storage has grown significantly, shifting markedly from analog to digital after 2000

Global installed, optimally compressed, storage



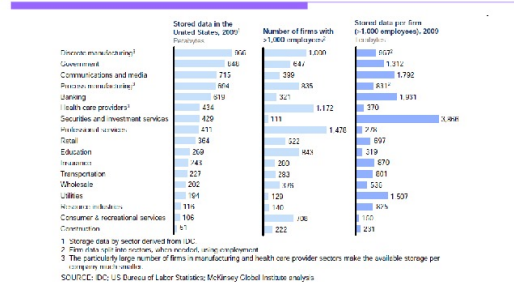
dalloz@dalloz.com

Global information created and available storage Exabytes



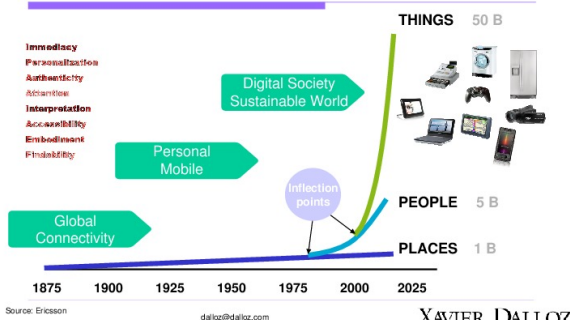
Source: IDC

Le tas d'or sur lequel les entreprises sont assises...



dalloz@dalloz.com

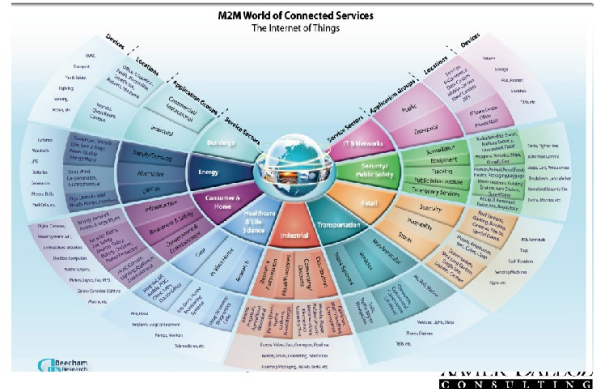
50 milliards d'objets connectés en 2025 des nouvelles créations de valeur



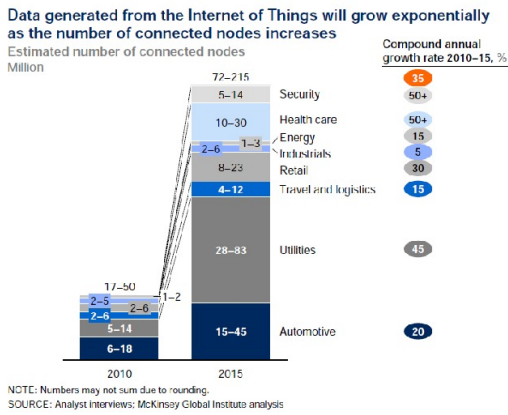
Source: Ericsson

dalloz@dalloz.com

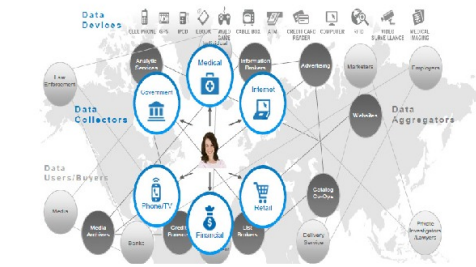
Nous n'avons encore rien vu !!



Source: M2M World of Connected Services



L'interaction avec notre environnement
 Vers une économie de l'immatériel et de l'intelligence collaborative



daloz@dalloz.com

XAVIER DALLOZ
 CONSULTING

« Le job en vogue dans les 10 prochaines années sera statisticien... La faculté de prendre des données— et être capable de les comprendre, de les traiter, d'en extraire de l'information, ainsi que de les visualiser et de les communiquer »

Hal Varian
 (Google's chief economist)

Nouvelles technologies, nouveaux métiers, nouveaux enjeux...

- De nouveaux enjeux
 - Comment exploiter ces nouveaux volumes de données ?
 - Comment les stocker?
 - Comment les traiter?
 - Comment les visualiser?
- Et de nouvelles technologies
 - Bases de données distribuées
 - Traitement de données distribué
 - Analyse d'événements en temps réel
 - "Cloud Computing"
 - Bases de données distribuées
 - Traitement de données distribué
 - Analyse d'événements en temps réel
 - "Cloud Computing"

daloz@dalloz.com

XAVIER DALLOZ
 CONSULTING

Exemples d'outils

- A/B testing
- Association rule learning
- Classification
- Cluster analysis
- Crowdsourcing
- Data fusion and data integration
- Data mining
- Ensemble learning
- Genetic algorithms
- Machine learning
- Natural language processing
- Neural Networks
- Network analysis
- Pattern recognition
- Predictive modeling
- Regression
- Sentiment analysis
- Signal processing
- Spatial analysis
- Statistics
- Supervised learning
- Simulation
- Time series analysis
- Unsupervised learning
- Visualization

daloz@dalloz.com

XAVIER DALLOZ
 CONSULTING

Le plus de données ouvertes, le plus d'intelligence à construire...

- Big Table
- Business intelligence
- Cassandra
- Data mart
- Data warehouse
- Hadoop
- MapReduce
- Metadata
- R.

daloz@dalloz.com

XAVIER DALLOZ
 CONSULTING

L'enjeu : créer de la valeur

- Il faut exploiter les **Big Data** non pour automatiser le passé,
- Mais pour faire:
 - **Plus** (productivité, pénétration...)
 - **Mieux** (compétitivité, pertinence...)
 - **Autre chose** (nouveaux marchés, innovations, partenariats, offres contextuelles...)
 - **Différemment**

daloz@dalloz.com

XAVIER DALLOZ
 CONSULTING

L'or des poussières

- Les techniques du **Big Data** facilitent des interactions
 - Moins chères
 - Plus rapides
 - Plus efficaces
 - Moins contraintes par le temps et les distances.
 - D'où de nouveaux modèles économiques
- De **plus en plus de stratégies gagnantes sont basées sur la valorisation**
 - de la poussière d'information (Google),
 - de petites transactions (eBay),
 - de petites publicités (Google encore),
 - de collaborations (logiciels libres)...
 - Appliquer à l'envers la règle du 80/20 pour transformer la poussière en or...

daloz@dalloz.com

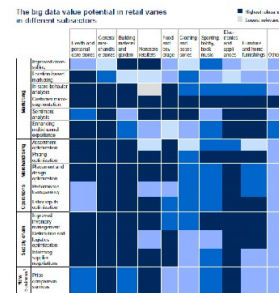
XAVIER DALLOZ
 CONSULTING

Exemples d'usage des Big Data

daloz@dalloz.com

XAVIER DALLOZ
 CONSULTING

Toutes les secteurs d'activité économique sont concernées...



daloz@dalloz.com

XAVIER DALLOZ
 CONSULTING

Les nouveaux acteurs utilisent massivement les Big Data

- **Zynga**
 - Créateur des jeux sociaux les plus populaires 10% de la population internet mondiale a joué un jeu Zynga (230M joueurs / mois) Déplace 1Pb de données chaque jour
 - Ajoute 1000 serveurs par semaine
- **Facebook**
 - 500M d'utilisateurs
 - 3.5B morceaux de contenu / semaine
 - 4B messages / jour
 - 1.2M photos / second (lues)
- **Twitter**
 - 70M Tweet / day
 - $800 \text{ tweets/sec} * 200b = 160kb/sec = 9Mb/min = 12Gb/jour$ de texte 8T de données / day 6B api calls / day

daloz@dalloz.com

XAVIER DALLOZ
 CONSULTING

Le commerce de détail

- La base de données de WalMart contient 2,5 petabytes :
 - CPFR
 - EDLP
 - Retail Link
- TESCO collecte 1,5 milliards de data chaque mois pour ajuster ses prix et les promotions
- eBay a une base de données de 6,5 et un datawarehouse de 2,5 petabyte

daloz@dalloz.com

XAVIER DALLOZ
 CONSULTING

Big data retail levers can be grouped by function

Function	Big data lever
Marketing	<ul style="list-style-type: none"> • Cross-selling • Location based marketing • In-store behavior analysis • Customer micro-segmentation • Sentiment analysis • Enhancing the multichannel consumer experience
Merchandising	<ul style="list-style-type: none"> • Assortment optimization • Pricing optimization • Placement and design optimization
Operations	<ul style="list-style-type: none"> • Performance transparency • Labor inputs optimization
Supply chain	<ul style="list-style-type: none"> • Inventory management • Distribution and logistics optimization • Informing supplier negotiations
New business models	<ul style="list-style-type: none"> • Price comparison services • Web-based markets

SOURCE: McKinsey Global Institute analysis

CONSULTING

La production des biens manufacturés (vendre avant de produire)

daloz@daloz.com

XAVIER DALLOZ CONSULTING

L'indicateur de base

$$\frac{\text{Résultat d'exploitation}}{\text{Capitaux fixe + circulant}} = \frac{\text{Résultat d'exploitation}}{\text{Chiffre d'affaire}} \times \frac{\text{Chiffre d'affaire}}{\text{Capitaux fixe + circulant}} = \text{Chiffre d'affaire} \times \text{Taux de rotation du capital d'exploitation}$$

Rentabilité économique
Elle croît si:

- ☐ Qualité, différence de l'offre,
- ☐ Efficacité de la relation client,
- ☐ Offre nouvelle,
- ☐ Nouvelles clientèles, extension territoriale

Augmente aussi si:
Efficacité accrue de l'utilisation du capital

Taux de marge
Il croît si l'on augmente

- ☐ la rentabilité des ventes
- ☐ leur pertinence (sur mesure)
- ☐ leur différenciation
- ☐ Leur valeur ajoutée

Taux de rotation du capital d'exploitation
la rentabilité financière croît avec

- ☐ la vitesse de rotation des actifs
- ☐ La réduction du capital fixe nécessaire (partenariats, entreprise étendue...)

daloz@daloz.com

XAVIER DALLOZ CONSULTING

3 impératifs liés par les Big Data: gérer, produire, vendre

- ☐ **Gérer: mieux exploiter le capital**
 - Le capital fixe
 - Le capital circulant
 - L'intelligence collective interne
 - Les synergies de réseau à valeur ajoutée
 - Les risques et opportunités (innovations, acquisitions...)
- ☐ **Produire: réduire les coûts**
 - Des matériaux, composants, équipements,
 - Personnels, management, frais généraux...
 - Coûts d'interactions
- ☐ **Vendre: accroître les revenus**
 - Parts de marché, nouveaux marchés
 - Clients: fidélisation, acquisition et extension
 - Time to market, cycle de renouvellement

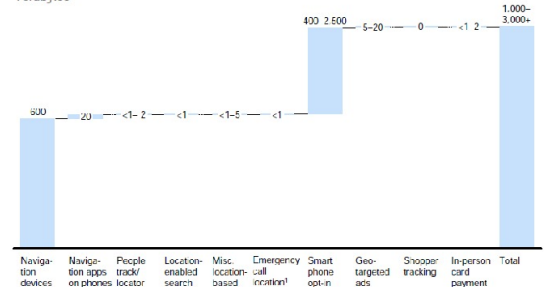
daloz@daloz.com

XAVIER DALLOZ CONSULTING

La géolocalisation

Overall personal location data total more than 1 petabyte per year globally

Personal location data generated globally, 2009



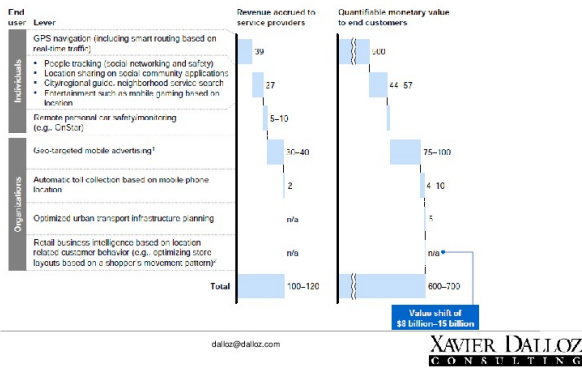
daloz@daloz.com

XAVIER DALLOZ CONSULTING

daloz@daloz.com

XAVIER DALLOZ CONSULTING

The value of the major levers increases to more than \$800 billion by 2020
\$ billion per annum

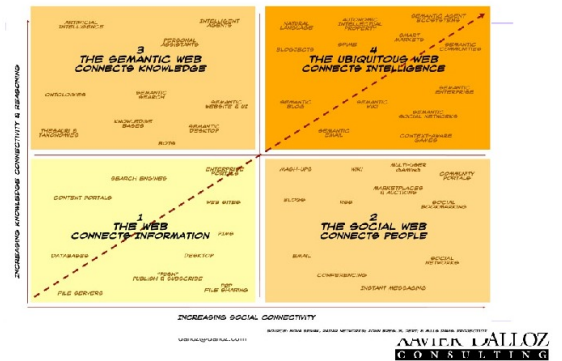


L'Internet des objets et la shazamisation

- Tags RFID
 - 12 millions vendus en 2011
 - 209 milliards en 2021
- De nouveaux process métiers (traçabilité)
 - Computer aided engineering
 - Computer aides manufacturing
 - Collaborative product development management
 - Digital manufacturing
 - Etc.

Le web sémantique

What is the evolution of the internet to 2020?



Recommandations

Definitions	Deep analytical	Big data savvy	Supporting technology
	People who have advanced training in statistics and/or machine learning and conduct data analysis	People who have basic knowledge of statistics and/or machine learning and define key questions, data can answer	People who serve as database administrators and programmers
Occupations ¹	<ul style="list-style-type: none"> • Actuaries • Mathematicians • Operations research analysts • Statisticians • Mathematical scientists • Industrial engineers • Epidemiologists • Economists 	<ul style="list-style-type: none"> • Business and financial managers • Biologists and medical research scientists • Engineers • Life scientists • Market research analysts • Survey researchers • Industrial organizational psychologists • Sociologists 	<ul style="list-style-type: none"> • Computer and information systems managers • Computer programmers • Computer software engineers for applications • Computer software engineers for system software • Computer system analysts • Database administrators

¹ Occupations are defined by the Standard Occupational Code (SOC) of the US Bureau of Labor Statistics and used as the proxy for types of talent in labor force. SOURCE: ICF Global and Xavier Dalloz, McKinsey Global Institute analysis.

daloz@dalloz.com XAVIER DALLOZ CONSULTING

Devenir problématique...

- Il faut apprendre à poser la bonne question
- Trop d'informations... tue l'information
- Réduire la complexité
- Mesurer pour valider les intuitions
- Les principaux gisements de création de valeurs sont :
 - "Predictive analysis"
 - Rendre l'information intelligible
 - Visualisation de données

2.2 Vincent HEUSCHLING (D-FI)

Cartographie des solutions *big data* du marché

Face au déluge de donnée, que nous vivons aujourd'hui, quelles sont les réponses des grands acteurs du marché ? Quels défis présente cette explosion du volume de données pour les infrastructures. Décryptage des offres des grands constructeurs que sont EMC, IBM, Oracle, ..., et des architectures innovantes du monde opensource.

Cartographie des solutions BigData

Panorama du marché et prospective

1

11 juin 2011

Solutions BigData

- Défi(s) pour les fournisseurs
- Quel marché
- Architectures
- Acteurs commerciaux
- Solutions alternatives

Solutions BigData le 9/6/2011 2

11 juin 2011 Vincent Heuschling

Quels Défis ?

- des volumes impossibles à traiter :
 - 30 To de logs par jour chez Facebook
 - 15 Po de data par an au CERN
- des croissances vertigineuses
- du business en temps réel
- des données différentes :
 - Non structurées, réparties, NoSQL...

Solutions BigData le 9/6/2011 3

11 juin 2011 Vincent Heuschling

Le quadrant magique (DW database management)

Source: Gartner (January 2011)

Solutions BigData le 9/6/2011 4

11 juin 2011 Vincent Heuschling

Positionnement des acteurs du marché

Solutions BigData le 9/6/2011 5

11 juin 2011 Vincent Heuschling

ROI

- Révolutionne les datawarehouses existants
- ROI de 27 mois à 6 mois
- 3 fois moins cher
- 4 fois plus rapide à implémenter

Solutions BigData le 9/6/2011 6

11 juin 2011 Vincent Heuschling

Architecture & composants

- Shared Disk vs Share Nothing Arch.
- Hadoop / HBase / HDFS
- Map Reduce

Solutions BigData le 9/6/2011 7 Vincent Heuschling
samed 11 juin 2011

Map Reduce

How much wood would a woodchuck chuck if a woodchuck could chuck wood

↓ Map Function :
output (word : 1)

how : 1
 much : 1
 wood : 1
 would : 1
 woodchuck : 1
 chuck : 1
 woodchuck : 1
 could : 1
 chuck : 1
 wood : 1

→

wood : 2
 woodchuck : 2
 chuck : 2
 : 1

Reduce Function :
output (word : sum(1))

Solutions BigData le 9/6/2011 8 Vincent Heuschling
samed 11 juin 2011

Map Reduce

- S'appuie sur une base key / value
- est scalable sur n serveurs
- permet d'enchaîner plusieurs Reduce
- beaucoup d'implémentations

Solutions BigData le 9/6/2011 9 Vincent Heuschling
samed 11 juin 2011

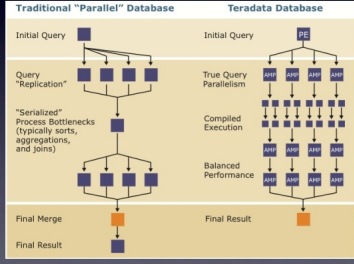
Acteurs du marché

- **Teradata**
- **Oracle / Exadata**
- **IBM / Netezza**
- **EMC / Greenplum**
- ...

Solutions BigData le 9/6/2011 10 Vincent Heuschling
samed 11 juin 2011

Teradata


- Depuis 1979
- Appliances
- Share nothing arch.
- Parallélisme
- Pour les DW
- De 6 To à 92 Po



Solutions BigData le 9/6/2011 11 Vincent Heuschling
samed 11 juin 2011

Oracle Exadata

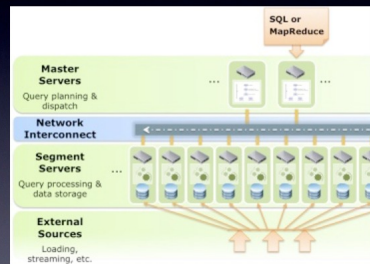
- «Database machine» (n'est pas une appliance)
- Serveurs de stockage (168 cores, 5 TB de flash cache, 45 TB utiles)
- Serveurs de traitements (128 cores / 2 TB de mémoire)
- 1500000 IOPS
- Data Load Rate: Up to 12 TB/hour



Solutions BigData le 9/6/2011 12 Vincent Heuschling
samed 11 juin 2011

EMC Greenplum

- Serveurs std
- Share nothing arch
- Map Reduce
- SQL



Solutions BigData le 9/6/2011

13

Vincent Heuschling

samed 11 juin 2011

13

IBM Netezza

- Blades IBM + Disques + FPGAs
- Share nothing arch.
- Map Reduce & SQL
- Data load rates de 2TB/h
- Produits : Skimmer (1TB à 10TB) & TwinFin (1TB à 1PB+)



Solutions BigData le 9/6/2011

14

Vincent Heuschling

samed 11 juin 2011

14

Alternatives et Opensource

- **Active circle**
- **Bases NOSQL**
- **Apache HADOOP**
- **Database.com**
- **Amazon Elastic Map Reduce**

Solutions BigData le 9/6/2011

15

Vincent Heuschling

samed 11 juin 2011

15

Active Circle

- FileSystem distribué
- Accès par NAS ou API
- Virtualisation sur disque et bande
- Noeuds locaux ou distants
- Réplication
- Hiérarchisation

Solutions BigData le 9/6/2011

16

Vincent Heuschling

samed 11 juin 2011

16

NOSQL : Not Only SQL

- Cassandra
- Google's BigTable : HBase
- MongoDB (documents, JSON)
- CouchDB (documents, JSON)

Solutions BigData le 9/6/2011

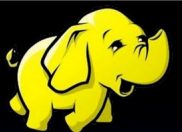
17

Vincent Heuschling

samed 11 juin 2011

17

Apache HADOOP



- HDFS (distributed high throughput FS)
- MapReduce
- HBase (scalable, distributed database)
- Hive (data warehouse infrastructure)
- Mahout (data mining library)
- Pig: (framework for parallel computation)
- ZooKeeper (distributed applications)

Solutions BigData le 9/6/2011

18

Vincent Heuschling

samed 11 juin 2011

18

Amazon Elastic Map Reduce

- Logique de PaaS :
 - Stockage avec Amazon S3
 - Processing avec un cluster Amazon EC2
- Mise en oeuvre instantanée
- Simple
- Economique (0,3 \$ / heure par node)

Solutions BigData le 9/6/2011

19

Vincent Heuschling

mercredi 11 juin 2011

19

Database.com

- Database as a Service (DaaS)
- Multi-tenant
- Scalable à l'infini
- économique : (\$10 / mois / 100000 records)

Solutions BigData le 9/6/2011

20

Vincent Heuschling

mercredi 11 juin 2011

20

Conclusions

- Des solutions dans la continuité de l'existant.
- Des innovations permettant des ROI attrayants : Attention aux ruptures
- Outils opensource en voie de maturation

Solutions BigData le 9/6/2011

21

Vincent Heuschling

mercredi 11 juin 2011

21

MERCI

Vincent Heuschling
vincent@heuschling.com
twitter : @vhe74

22

mercredi 11 juin 2011

22

2.3 Alexandru Costan (INRIA)

Analyse des systèmes de stockage à grande échelle pour les applications de traitement intensif des données

Avec l'augmentation rapide des volumes de données dans de nombreux domaines d'application de la science de l'ingénierie et des services d'information, les défis posés par les traitements intensifs des données présentent une importance croissante. Avec l'émergence des infrastructures récentes (plateformes de type *Cloud*, architectures massivement parallèles pétaflopiques) réaliser une gestion des données capable de passer à l'échelle dévient un enjeu crucial car les performances globales des applications dépendent des propriétés du service de gestion des données.

Nous définissons un ensemble de principes pour la conception de systèmes de stockage distribués, optimisés pour pouvoir passer à large échelle, tout en permettant autant de manipulations concurrentes des données que possible. Combinés, ces principes peuvent aider les développeurs de systèmes de stockage distribués à répondre aux exigences strictes de gestion de données à grande échelle.

Nous analysons ensuite plusieurs systèmes de stockage représentatifs afin d'évaluer la façon dont ils se conforment à ces principes de conception. Nous nous concentrons sur les systèmes de fichiers spécialisés qui ont été introduits pour cibler spécifiquement les besoins des applications de traitement intensif des données : HDFS, la couche de stockage standard utilisé par Hadoop MapReduce, GPFS proposé par IBM, ainsi que les systèmes de fichiers distribués massivement parallèles, comme Lustre ou PVFS, généralement utilisé pour les *clusters* de calcul à grande échelle. Avec l'émergence du calcul de type *Cloud*, des solutions de stockage spécialement conçues pour s'adapter à ce contexte ont été développés : nous présenterons textttAmazon S3.

Nous détaillons en particulier les avantages potentiellement importants du versionnage pour améliorer les performances d'accès hautement concurrents aux données des applications. Dans ce contexte, nous proposons une interface d'accès basé sur la gestion de versions des données, matérialisée au sein de la plate-forme `BlobSeer` développée par l'équipe `KerData` de l'INRIA à Rennes. Cette approche qui permet d'exploiter d'une manière efficace le parallélisme inhérent des flux de données : nous en illustrons utilisation avec une application d'analyse conjointe de données génétiques et de neuro-imagerie.

A Survey of Large Scale Storage Systems for Data Intensive Applications

Alexandru Costan
 KerData research team
 INRIA Rennes - Bretagne Atlantique, France

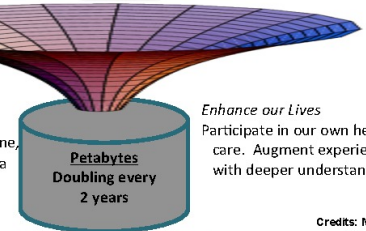
Séminaire Aristote, Ecole Polytechnique, 9 June 2011



Context Today: an Explosion of Data



The Challenge:
 Enable Discovery.
 Deliver the capability to mine, search and analyze this data in near real time.



Enhance our Lives
 Participate in our own health care. Augment experience with deeper understanding.

Credits: Microsoft



New Challenges for Large-scale Data Storage

Important issues:

- Scalable architecture (10⁶ nodes)
- Massive unstructured data (Terabytes)
- Many data objects (10⁹)
- Transparency
- High concurrency (10⁶ concurrent clients)
- Fine-grain access (Megabytes)

Applications: distributed, with high-throughput requirements under concurrency

- Map-Reduce-based data-mining applications
- Governmental and commercial statistics
- Data-intensive HPC simulations
- Checkpointing for massively parallel computations
- On-Line Social Networks

Target platforms: from large clusters, grids and desktop grids to clouds and petascale machines



Big Data storage systems design principles

Data organization

- Scalability, transparency
- Ex: distributed files systems, object based storage devices (OSDs)

Asynchronous management

- Atomicity

Concurrency control

- Application-level parallelism
- Ex: locks, versioning

Data striping

- Configurable chunk distribution strategy
- Dynamically adjustable chunk sizes

Distributed metadata management

- Data availability



Specialized distributed storage systems

Data-intensive oriented file systems

- GFS
- HDFS

Parallel file systems

- GPFS
- Lustre

Cloud data storage services

- S3
- Azure



Specialized distributed storage systems

Data-intensive oriented file systems

- GFS
- HDFS

Parallel file systems

- GPFS
- Lustre

Cloud data storage services

- S3
- Azure



HDFS (Hadoop Distributed File System)



Part of Yahoo! Hadoop

- MapReduce implementation
- Open Source
- Java based



Distributed storage system

- Files are divided into large blocks (64 MB)
- Blocks are distributed across the cluster
- Blocks are replicated to help against hardware failure
- Data placement is exposed so that computation can be migrated to data

Notable differences from mainstream DFS work

- Single 'storage + compute' cluster vs. separate clusters
- Simple I/O centric API

HDFS Architecture: NameNode (1)



Master-Slave Architecture

HDFS Master "NameNode"

- Manages all file system metadata in memory
 - List of files
 - For each file name, a set of blocks
 - For each block, a set of DataNodes
 - File attributes (creation time, replication factor)
- Controls read/write access to files
- Manages block replication
- Transaction log: register file creation, deletion, etc.

HDFS Architecture: DataNodes (2)



HDFS Slaves "DataNodes"

A DataNode is a block server

- Stores data in the local file system (e.g. ext3)
- Stores meta-data of a block (e.g. CRC)
- Serves data and meta-data to Clients

Block Report

- Periodically sends a report of all existing blocks to the NameNode

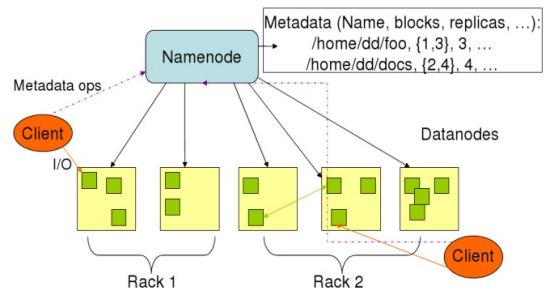
Pipelining of Data

- Forwards data to other specified DataNodes

Perform replication tasks upon instruction by NameNode

Rack-aware

HDFS Architecture (3)



Fault Tolerance in HDFS



DataNodes send heartbeats to the NameNode

- Once every 3 seconds

NameNode uses heartbeats to detect DataNode failures

- Chooses new DataNodes for new replicas
- Balances disk usage
- Balances communication traffic to DataNodes

Data corectness

- Use Checksums to validate data: CRC32

NameNode failures

- Single point of failures

Data-intensive oriented file systems



Huge files

Structured storage can be built on top

Fine grain concurrent reads

Pros

- No locking
- Data location aware

Cons

- Centralized metadata
- Expensive updates

Specialized distributed storage systems



Data-intensive oriented file systems

- GFS
- HDFS

Parallel file systems

- GPFS
- Lustre

Cloud data storage services

- S3
- Azure

GPFS (General Parallel File System)



Developed by IBM

- High-performance shared-disk clustered file system
- Used by many supercomputers in Top500

Distributed storage system

- Files are divided into *small blocks* (less than 1 MB)
- Blocks are distributed across the cluster
- Blocks are *RAID-replicated* or file system node replicated
- *Transparent* data location



Notable differences

- *Distributed metadata*
- Efficient indexing of directory entries for very large directories.
- *POSIX semantics*
- Network partition aware

GPFS Architecture - Special Node Roles



File system nodes

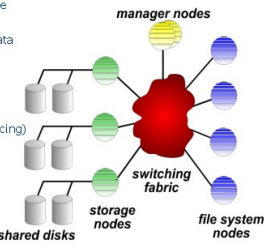
- Run user programs, read/write data to/from storage nodes
- Cooperate with manager nodes to perform metadata operations

Storage nodes

- Implement block I/O interface
- Interact with manager nodes for recovery (e.g. fencing)
- Data and metadata striped across multiple disks - multiple storage nodes

Manager nodes

- File system configuration: recovery, adding disks
- Disk space allocation manager, quota manager
- File metadata manager - maintains file metadata integrity
- Global lock manager



Credits: IBM

Lustre



Massively parallel distributed file system (owned by Oracle)



Used by most supercomputers:

- The world's fastest computer - Tianhe-1A
- Jaguar (ORNL), LBNL, CEA

Features:

- OSD based
- Open source

Lustre Architecture



Metadata Server (MDS)

- Active / Passive
- Filenames, directories, access permissions, file layout

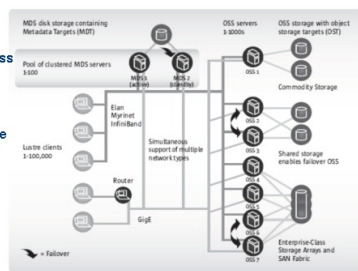
Object Storage Servers (OSS)

- Store data on Object Storage Targets
- Distributed locking

Clients

- POSIX semantics

Fault tolerance: "failure as an exception"



Parallel file systems



Mounted as regular file systems

Data striping

Advanced caching

Pros

- Distributed data
- MPI optimizations

Cons

- Locking-based
- Too many small files

Specialized distributed storage systems



Data-intensive oriented file systems

- GFS
- HDFS

Parallel file systems

- GPFS
- Lustre

Cloud data storage services

- S3
- Azure

S3



Amazon Simple Storage Service:

- "storage for the Internet"
- (cheap) pay per use policy (for storage, requests, data transfers)

Design

- Objects (up to 5TB) stored in *buckets*, identified using *keys*
- Buckets stored in one of several Regions
- Clients authorization using ACLs
- Access through Web interfaces: REST, SOAP, BitTorrent



Notable uses

- FUSE – allows EC2-hosted Xen images mount an S3 bucket as a file system
- Apache Hadoop
- Tumblr

Azure



Proposed by Microsoft within Windows Azure PaaS cloud

Data manipulation based on HTTP

All data replicated 3 times

Blobs

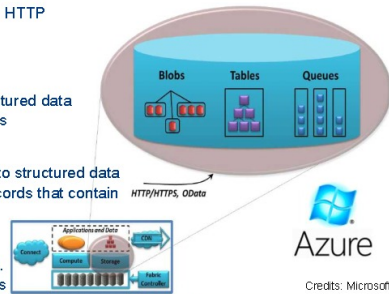
- Up to 1 TB of unstructured data
- Grouped in containers

Tables

- Fine grained access to structured data
- Group of entities / records that contain properties

Queues

- Asynchronous comm. between cloud instances



Cloud data storage services



Virtualize storage resources

Pay for duration, size and traffic

Flat naming scheme

Simple access model

Pros

- High data availability
- Versioning

Cons

- Limited object size
- Low throughput

Limitations of existing approaches



Issue	Parallel FS	Data-intensive FS	Cloud store
Too many small files	✗	✓	✗
Centralized metadata	✓	✗	✓
No versioning support	✗	✗	✓
No fine grain writes	✓	✗	✗

✓ = addressed issue

[Nicolae et al., 2010]

Concurrency-optimized BLOB management: The BlobSeer Approach



BlobSeer: software platform for scalable, distributed BLOB management

- Huge data (TB) - BLOBs: Binary Large Objects
- Highly concurrent, fine-grain access (MB): Read/Write/Append
- Developed by the KerData team at INRIA, Rennes

Overview of key design choices

- Decentralized data storage
- Decentralized metadata management
- Versioning-based concurrency control, multiversioning exposed to the user
- Lock-free concurrent writes (enabled by versioning)

A back-end for higher-level, sophisticated data management systems

- Short term: highly scalable distributed file systems
- Middle term: storage for cloud services
- Long term: extremely large distributed databases

<http://blobseer.gforge.inria.fr/>

BlobSeer: Key Design Choices



Distributed data

- Each BLOB is fragmented into "chunks" (pages)
- Huge data amounts to be distributed all over the storage nodes
- Reduced contention for simultaneous accesses to disjoint parts of the BLOB

Distributed Metadata

- Goal: locate chunks that make up a given BLOB
- Fine-grained and distributed
- Efficiently managed through a segment tree over a DHT

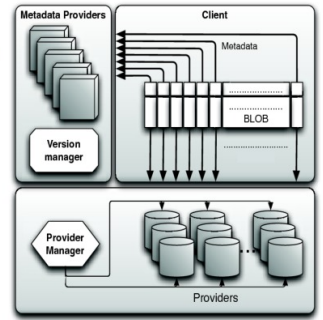
Versioning-based concurrency control

- Update/append: generate new chunks rather than overwrite
- Metadata is extended to incorporate the update
- Both the old and the new version of the BLOB are accessible
- Lock-free approach

<http://blobseer.gforge.inria.fr/>

BlobSeer: Architecture

- Clients**
 - Perform fine grain BLOB accesses
- Providers**
 - Store the chunks of the BLOB
- Provider manager**
 - Monitors the providers
 - Favors data load balancing
- Metadata providers**
 - Store information about chunk location
- Version manager**
 - Ensures concurrency control



Integrating BlobSeer in the Hadoop Map-Reduce Framework



MapReduce: a natural application class for BlobSeer:

- Case study: Yahoo!'s Hadoop MapReduce framework
- Approach: use BlobSeer instead of Yahoo!'s Hadoop file system (HDFS)
- Motivation: HDFS has limited support for concurrent access to shared data

Implementing the HDFS API for BlobSeer

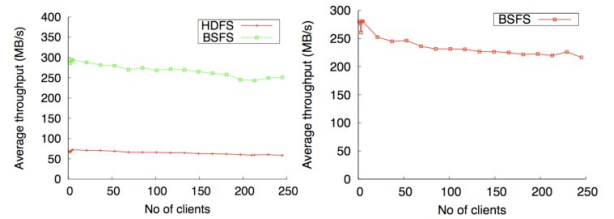
- Implements basic file system operations: create, read, write...
- Introduces support for concurrent append operations

BlobSeer File System (BSFS)

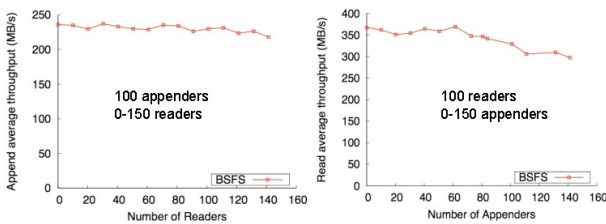
- File system namespace - keeps file metadata, maps files to BLOB's
- Client-side buffering: data prefetching, write aggregation
- Exposes data layout to Hadoop, just like HDFS

BSFS vs. HDFS

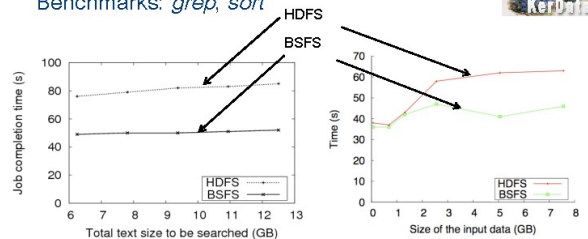
Concurrent Reads, Concurrent Appends



BlobSeer Scales Up: Readers and Writers Do Not Interfere (Almost!)



BSFS Does Better Than HDFS! Benchmarks: *grep*, *sort*



- Relevant publications
- JPDC (2010), Special Issue on Data-Intensive Computing
 - IEEE IPDPS 2010
 - MapReduce 2010 (held in conjunction with HPDC 2010)

The AzureBrain Project: BlobSeer on Microsoft Azure Clouds

- Application**
- Large-scale Joint Genetic and Neuroimaging Data Analysis
- Goal**
- Assess and understand the variability between individuals
- Approach**
- Optimized data processing on Microsoft's Azure clouds based on the BlobSeer concurrency-optimized platform
- INRIA teams involved**
- KerData (Rennes)
 - PARIENTAL (Saclay)
- Framework**
- Joint MSR-INRIA Research Center
 - MS involvement Azure teams, EMIC

Microsoft et l'INRIA annoncent un partenariat autour du Cloud Computing

Posté le 28/10/2010 | Institutionnel

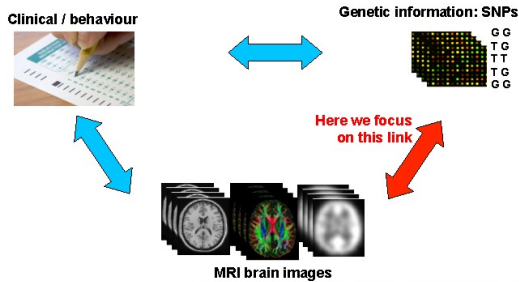
Pendant 2 ans, le projet AzureBrain, au service de la neuro-imagerie, va bénéficier des solutions Microsoft Windows Azure

Issy-les-Moulineaux, le 28 octobre 2010 – Microsoft et l'INRIA (Institut National de Recherche en Informatique et en Automatique) renforcent aujourd'hui leur collaboration avec le lancement du projet de recherche AzureBrain qui sera réalisé au sein du centre de recherche commun INRIA-Microsoft Research. AzureBrain a pour objectif de permettre des avancées précieuses dans le domaine de la neuroscience et de la neuro-imagerie. Microsoft met concrètement au service de ce projet, des ressources de Cloud Computing qui permettent d'accélérer le rythme des recherches, offrant des puissances de calcul et de traitement sur-mesure, à travers les datacenters.

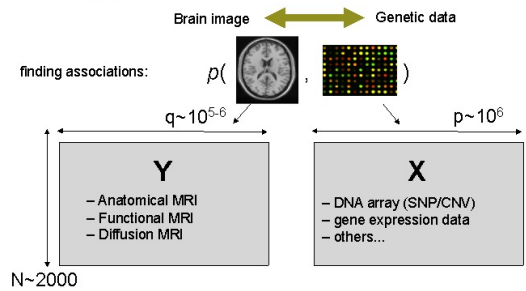
NeuroSpin Neuroimaging center at Saclay

- IRM 3T
- IRM 7T
- IRM 11.74T
- logitech
- Futur scanner 17.657060tesla
- Clinical area (3 hospital beds, neuro psy rooms, EEG / A/EG)
- Library and conferences
- Laboratoire and offices
- Pre-clinical area (transgenic mice, primates, open access fMRI, etc)

The Imaging Genetics Challenge: Comparing Heterogeneous Information

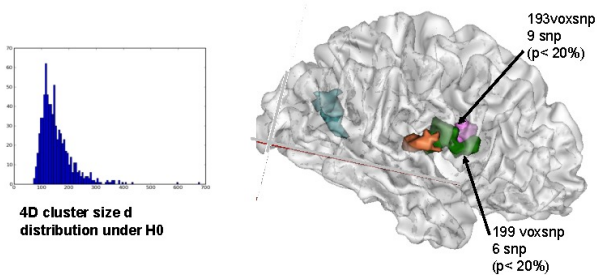


Imaging Genetics Methodological Issues



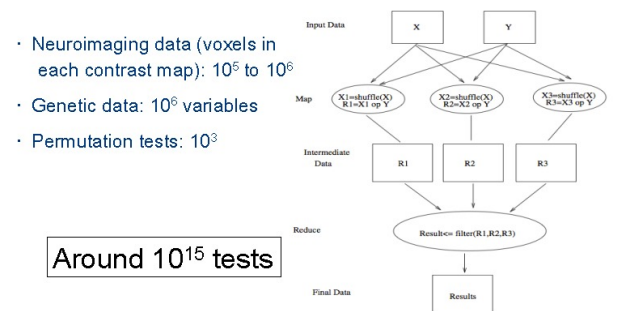
Approach: Searching Statistical Associations Between Pairs

Illustration



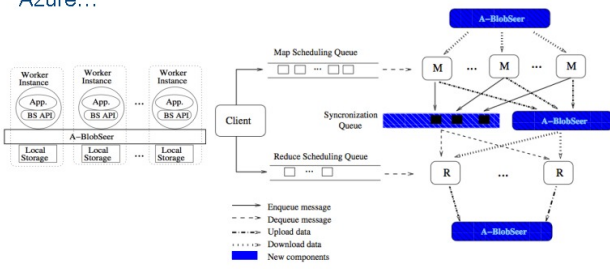
A-Brain:
The goal is to reproduce this kind of study with 10^5 larger data

The Computational Problem

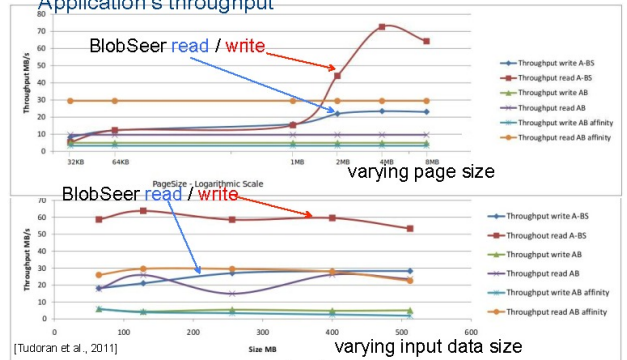


BlobSeer as a storage backend in Azure...

... used within the MapReduce based application



Application's throughput



Summary

Difficult to maximize all the objectives: achieve a very **high data throughput** for **highly concurrent, fine-grain data accesses**

Concurrency control based on **locking mechanisms** often creates bottlenecks

Object based storage approaches ensure scalability

Consistency model: **CAP**

Data-intensive specific solutions exploit application level parallelism but force users to adhere to a **specific programming paradigm**

Thank you!



For more information...

- **BlobSeer**: <http://blobseer.gforge.inria.fr>
- **KerData**: <http://iris.a.fr/kerdata>

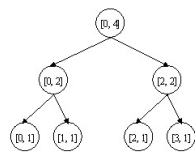
Metadata Zoom (1)

Organized as a segment tree

Each node covers a range of the blob identified by (offset, size)

The first/second half of the range is covered by the left/right child

Each leaf corresponds to a chunk and holds information about its location



Metadata Zoom (2)

Each node holds versioning Information

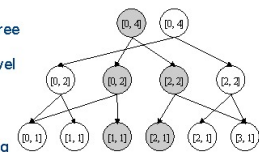
Write/Append

- Add leaves and build subtree up to the root
- The tree may grow one level

Read: descend from the root towards the leaves

Tree nodes are **distributed** among metadata providers

Highly scalable access concurrency: **R/R, R/W, W/W**



Metadata Zoom (2)



Each node holds versioning Information

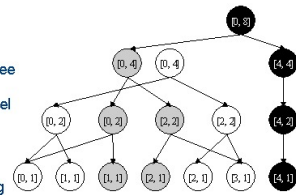
Write/Append

- Add leaves and build subtree up to the root
- The tree may grow one level

Read: descend from the root towards the leaves

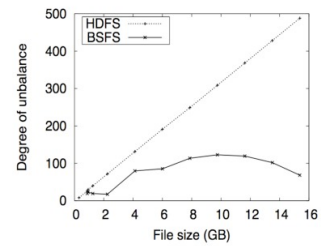
Tree nodes are **distributed** among metadata providers

Highly scalable access
concurrency: **R/R, R/W, W/W**

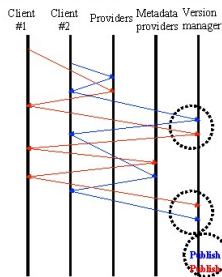


BSFS vs. HDFS

Load balancing the storage nodes



How Versioning Enables Efficient, Heavy Access Concurrency



Chunks are written concurrently by the clients

Then, versions are assigned in the order the clients finish writing

Then, metadata is written concurrently by the clients

Versions are published in the order they were assigned

Leveraging BlobSeer on Clouds: MapReduce



MapReduce: a **simple programming model** for data-intensive computing

Typical problem solved by MapReduce

- Read a lot of data
- **Map**: extract something you care about from each record
- Shuffle and Sort
- **Reduce**: aggregate, summarize, filter, or transform
- Write the results

Approach: **hide messy details** in a runtime library

- Automatic parallelization
- Load balancing
- Network and disk transfer optimization
- Transparent handling of machine failures

Implementations: Google MapReduce, Hadoop (Yahoo!)

2.4 Joerg Bienert (ParStream)

An innovation solution to manage heterogeneous big data

Analyzing Big Data is getting more and more important for companies in all industries, *e.g.* web analytics, fraud detection, smart metering, telco, *etc.* Current established database products are not able and not designed to perform large scale data analytics. New approaches like Map/Reduce lack important features like short response times. `ParStream` is a novel innovative database technology focusing on processing large data sets (billions of records) in milliseconds and with low latency. `Parstream` is a columnar in memory database running on multiprocessor architectures and, as first product, on GPU based HPC-Servers.

ParStream
Big Data – Real Time – Low Latency
Ecole Polytechnique
June 9th 2011

Agenda

- Big Data – The Challenge
- Current Approaches
- ParStream – A HPC Database for Big Data
- Use Cases

2

Big Data

The amount of digital information increases tenfold every five years



- "We are at a different period because of so much information", says James Cortada of IBM
- "...what is scarce is the ability to extract wisdom from them.", Hal Varian, Google's chief economist
- "Every day I wake up and ask, how can I flow data better, manage data better, analyse data better?", says Rollin Ford, the CIO of Wal-Mart.
- "The data-centered economy is just nascent", admits Mr. Mundie of Microsoft.

3

Big Data

„Analyzing big data—will become a key basis of competition, underpinning new waves of productivity growth, innovation, and consumer surplus“



- Making big data more accessible in a timely manner
- Using data and experimentation to expose variability and improve performance
- Segmenting populations to customize actions
- Replacing and supporting human decision-making with automated algorithms
- Innovating new business models, products, and services.

4

Challenges in Big Data Analytics

There are challenges in Big Data Analysis in different dimensions

Billions of Records 	Thousands of Columns 	Concurrent Queries 	Real Time & Low Latency 	Complex Queries
-------------------------	--------------------------	------------------------	-----------------------------	---------------------

5

Big Data

eCommerce & Web Analytics

A collection of logos for major e-commerce and web analytics companies: eBay, Amazon, WebTrends, Google Analytics, OMNITURE, and salesforce.com. A small screenshot of a data dashboard is also visible on the right.

6

Big Data 
 Social Media Analytics



7

Big Data in Telecommunications 
 Billing, Pricing, Fraud detection



8

Big Data in Utilities 
 Smart metering & Smart grid



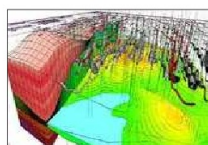
9

Big Data in Finance 
 Share Prices History, Fraud Detection, Algo-Trading



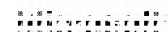
10

Big Data in other Industries 
 Retail, Telematics, Science, Mining



11

Established Databases 
 Current databases are not engineered for mass-data




“Existing Database Architectures are 20-30 years old and are not able to cope with current data sizes.”

Several Statements on Conference for Very Large Databases VLDB 2010, Singapore



12

Hadoop Ecosystem 

Some New approaches not suitable for Realtime Processing

"MapReduce isn't suited to calculations that need to occur in near real-time. You can't do anything that takes a relatively short amount of time, so we got rid of it."

Lipkowitz, Senior Director of Engineering, Google

13

Big Data Threat 

Big Data Analysis requires new "engines"



14

Big Data products 

New Marketplayers – and their downsides

 VERTICA No Partitioning	 Greenplum No Index, Postgres based
 aster data Map Reduce approach	 NETEZZA Hardware bound
 EXASOL No Index	 PAR)ACCEL No Index

15

Big Data 

New Market Players – and their setbacks

 hp No Partitioning	 VERTICA No Index, Postgres based
 TERADATA Map Reduce approach	 IBM Hardware bound
 aster data No Index	 NETEZZA No Index
 EXASOL No Index	 PAR)ACCEL No Index


16

ParStream 





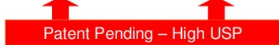


The first HPC Database focusing on the Real-Time Analysis of Big Data

17

ParStream – Building Blocks 

ParStream combines state of the art database technologies with unique technologies

 Column Store	 In-Memory Technology	 High Performance Index	 Custom Query Operators
Fast data access for analytical processing	Data and indices can stay in memory due to efficient compression	Unique index structure allows highly parallel execution of queries	SQL, JDBC and powerful C++ API enable fast query processing
 Patent Pending – High USP			

18

High Performance Indexing

Fast and flexible analytics through highly compressed and partitioned Bitmap-Indexes

The diagram illustrates a data processing pipeline: **Big Data** (represented by a stack of cylinders) is processed into an **Index** (represented by a folder icon). This index is then subjected to **Compression** (represented by a folder icon with a compression symbol) and **Partition** (represented by a folder icon with a partition symbol). Finally, the data is executed in **Parallel Execution** (represented by multiple server rack icons).

19

High Level Technical Architecture

ParStream is a complete database on a highly optimized architecture

The architecture diagram shows **Data Source** feeding into **ETL** processes. These feed into a **Loader** within the **ParStream** database. The ParStream database consists of a **Row-oriented record store**, a **Column-oriented record store**, and an **Index Store**. These are connected to **Query Logic**, which outputs **Results** in response to a **Query**. The entire system runs on **Single Servers / Clusters / Clouds with CPUs & GPUs**. Key optimization features include **High Speed Import & Index**, **Compression Partitioning**, and **Caching Parallelization**.

20

ParStream - Product Features

ParStream is optimized for all challenges in Big Data analysis

Low Latency	Clustering	Scalability	Editions
Queries include new data immediately after continuous import	Configurable partitioning allows horizontal and vertical clustering (fault tolerant)	ParStream scales linearly up to PetaByte	SW Package for Linux on Standard HW (Windows upc.) GPU Server Appliance

21

GPU Processing – motivation

Graphical Processing Units (GPUs) outperform standard processors

The Performance Gap Widens Further

The graphs show a significant performance gap between GPUs and CPUs. The left graph shows **Peak Single Precision Performance (GFlops/sec)** for NVIDIA GPUs (Tesla 8-series, 10-series, 20-series) and X86 CPUs (Nehalem 3 GHz). The right graph shows **Peak Memory Bandwidth (GB/sec)** for the same series. Both show exponential growth for GPUs compared to the linear growth of CPUs.

22

Cloud Ready

ParStream delivers outstanding performance on all infrastructure Setups

The image shows a server rack on the left and the ParStream logo inside a cloud on the right, indicating its cloud-ready nature.

23

Performance ParStream vs PostgreSQL

ParStream delivers in sub-seconds even on > 100 million rows
ParStream clearly dominates PostgreSQL and others

The left graph, titled **PostgreSQL (blue) reaching its limits**, shows response time increasing sharply as the number of rows increases. The right graph, titled **ParStream delivering in sub-second**, shows response time remaining consistently low (under 1 second) even for up to 180 million rows.

24

Explorative analytics on billions of records

ParStream analyzes billions of records within milliseconds 




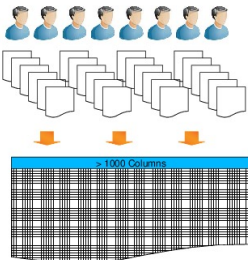
Case: Web Analytics

- Unique Visitor Calculation (select count distinct)
- Analysis of 8 billion records in 15 ms on 4-server cluster

25

Flexible analytics on thousands of columns

ParStream allows querying data by selecting any combination of different columns 





Case: Market Research

- Pattern recognition
- Flexible multi-column filtering & grouping
- 20 million records with 1000 data columns
- 5000 queries in 5 seconds

26

Concurrent queries and high throughput

ParStream can execute many queries efficiently at the same high throughput 



Case: Online Travel Search

- High throughput even with complex filter & sorting criteria
- 100 queries per second per node on 1 billion travel offers
- Throughput scales linearly with number of nodes

27

Use case: travel search engine

Comparison for a search appliance delivering 100 travel offers out of 1 billion based on 25 independent, optional filter criteria

Response Time

Database X: 6.5 sec

Parstream: 0.08 sec

factor 80

Memory Requirements

Database X: 25 Index size, ~5000MB for each query

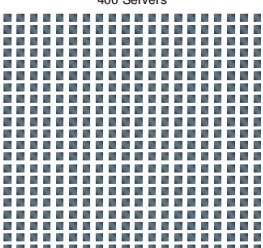
Parstream: 24, ~5MB for each query

factor 1000

28


Use case: travel search engine

Server Requirements




Sizing for a search appliance executing 1000 queries per second to deliver 100 travel offers out of 1 billion based on 25 optional filter criteria

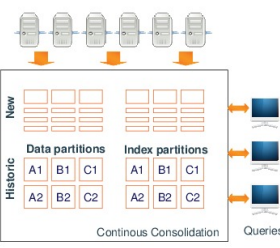
factor 20



29

Instant analytics on up-to-the-second data


Simultaneous data import and querying allows comparison of current & historical data 





Case: Web Analytics

- Continuous import of newly created data every second
- Simultaneous import and querying
- Single node performance typically 1TB per hour
- Continuous trend spotting for alert-function

30

Advanced analytics 


Standard SQL features and custom query nodes provide maximum flexibility and speed 



Case: Climate Research

- Identification of hurricane risk areas
- Geo-clustering with customized query node
- Implementation of custom algorithms via C++ API
- Analysis of 3 billion records in 100 ms

31

Testimonials 

Leading experts, customers and hardware providers have realized the benefits of ParStream

Experts are impressed...

- "An extremely innovative idea"
- "Convinced me completely"
- "2 years ahead of competition"

*Prof. Dr. Markl,
leading database expert at TU Berlin,
Formerly at IBM for DB2 optimization*


Customers are satisfied...

- Outperforms all databases by a factor of at least 35
- Reduced our response time by factors up to 12.000
- Scales linearly, runs stable in production


Leading web-analytics company

Awards...

- "One to watch award 2010" from Nvidia, Cooley & Adobe
- "First database running on high-performance GPUs"




32

Summary 



- Analyzing Big Data is getting more and more important for every industry
- Existing Database Architectures are 25-30 years old and not designed for Big Data
- New approaches often do not fulfill promises
- ParStream is a new powerful database, specifically designed for Big Data Analysis

33

Contact 

Questions ?

joerg.bienert@parstream.com
Phone +49 221 97 76 14 80

Geflügel-AG ist
Rechenzentrum für Statistik
und Marketing

Informationen über unsere
Anforderungen und Services

34

2.5 Ronan Keryell (HPC-Project)

Environnement de programmation pour traitements massifs sur architectures modernes

L'ubiquité de l'informatique déclenche une avalanche de données à traiter de manière rapide et économique. Malheureusement, pour des raisons physiques, la loi de Moore ne fournit plus de processeurs plus rapides (dissipation thermique, vitesse de la lumière, ...) mais fournit néanmoins toujours plus de transistors. Le seul moyen d'utiliser ces transistors est d'utiliser du parallélisme massif, mais cela remet en question les architectures et les modes de programmation. Les architectures actuelles (GPU, MP-SOC, Tiler, FPGA, ...) seront présentées avec leurs avantages et leurs contraintes ainsi que leurs pendants logiciels permettant de les utiliser au mieux. L'environnement de compilation source à source `Par4All` de `HPC Project` est un moyen de s'abstraire de certains détails de programmation.

Programming environments for big data processing on modern parallel architectures

Ronan KERYELL
HPC Project
2011/06/09

Le déluge de données, comment en tirer parti ?
Séminaire Aristote


HyperParallel Technologies (1992–1998)

- Parallel computer
- Proprietary 3D-torus network
- DEC Alpha 21064 + FPGA
- HyperC (follow-up of PompC @ LI/ENS Ulm)
 - ▶ PGAS (Partitioned Global Address Space) language
 - ▶ An ancestor of UPC...
- Already on the Saclay Plateau ! ☺

Quite simple business model

- Customers need just to rewrite all their code in HyperC ☺
- Difficult entry cost... ☹

- Niche market... ☹
- American subsidiary with dataparallel datamining application acquired by Yahoo! in 1998
- Closed technology → lost for customers and... founders ☹




HyperParallel Technologies (1992–1998)

- Parallel computer
- Proprietary 3D-torus network
- DEC Alpha 21064 + FPGA
- HyperC (follow-up of PompC @ LI/ENS Ulm)
 - ▶ PGAS (Partitioned Global Address Space) language
 - ▶ An ancestor of UPC...
- Already on the Saclay Plateau ! ☺

Quite simple business model

- Customers need just to rewrite all their code in HyperC ☺
- Difficult entry cost... ☹

- Niche market... ☹
- American subsidiary with dataparallel datamining application acquired by Yahoo! in 1998
- Closed technology → lost for customers and... founders ☹




Present motivations: reinterpreting Moore's law (I)

The good news ☺


- Number of transistors still increasing
- Memory storage increasing (DRAM, FLASH...)
- Hard disk storage increasing
- Processors (with captors) everywhere
- Network is increasing

- The bad news ☹
 - ▶ Transistors are so small they leak... Static consumption
 - ▶ Superscalar and cache are less efficient compared to transistor budget
 - ▶ Storing and moving information is expensive, computing is cheap: change in algorithms...
 - ▶ Light's speed as not improved for a while... Hard to reduce latency
 - Chips are too big to be globally synchronous at multi GHz ☹
 - ▶ pJ and physics become very fashionable



Present motivations: reinterpreting Moore's law (II)

- ▶ Power efficiency in $\mathcal{O}(\frac{1}{f})$
 - Transistors cannot be used at full speed without melting ☹
 - Heat
- ▶ I/O and pin counts
 - Huge time and energy cost to move information outside the chip ☹
- Rotating hard disk with 1D density d increase
 - ▶ Storage in $\mathcal{O}(d^2)$
 - ▶ But track speed only $\mathcal{O}(d)$
 - Reading all the disk in $\mathcal{O}(\frac{1}{f})$ ☹




Heterogeneous parallelism

Parallelism is the only way to go...

- Research is just crossing reality!
- Scaring... ☹
- Exciting! ☺

No one size fit all...

Future will be heterogeneous




The "Software Crisis"

Edsger DIJKSTRA, 1972 Turing Award Lecture, « The Humble Programmer »

"To put it quite bluntly: as long as there were no machines, programming was no problem at all; when we had a few weak computers, programming became a mild problem, and now we have gigantic computers, programming has become an equally gigantic problem."


http://en.wikipedia.org/wiki/Software_crisis
 ⚠ But... it was before parallelism democratization! ☺



Programming environments for big data processing on modern parallel architectures
 HPC Project Ronan KERYELL 6 / 101

Time to be back in parallelism!

- Good time for more start-ups! ☺
- ~ 2006: thinking to yet another start-up...
- People that met ≈ 1990 at the French Parallel Computing military lab SEH/ETCA
- Later became researchers in Computer Science, CINES director and ex-CEA/DAM, venture capital and more: ex-CEO of Thales Computer, HP marketing...
- HPC Project launched in December 2007
- ≈ 30 colleagues in France (Montpellier, Meudon), Canada (Montréal with Parallel Geometry) & USA (Mountain View)




Programming environments for big data processing on modern parallel architectures
 HPC Project Ronan KERYELL 7 / 101


HPC Project hardware: WildNode from Wild Systems

Through its Wild Systems subsidiary company

- WildNode hardware desktop accelerator
 - ▶ Low noise for in-office operation
 - ▶ x86 manycore
 - ▶ nVidia Tesla GPU Computing
 - ▶ Linux & Windows



<http://www.wild-systems.com>




Programming environments for big data processing on modern parallel architectures
 HPC Project Ronan KERYELL 8 / 101

HPC Project software and services

- Parallelize and optimize customer applications, co-branded as a bundle product in a WildNode (e.g. Presagis Stage battle-field simulator, WildCruncher for Scilab/...)
 - ▶ Acceleration software for the WildNode
 - ▶ GPU-accelerated libraries for C/Fortran/Scilab/Matlab/Octave/R
 - ▶ Transparent execution on the WildNode
- Remote display software for Windows on the WildNode

HPC consulting


- Optimization and parallelization of applications
- *High Performance?*... not only TOP500-class systems: power-efficiency, embedded systems, green computing...
- ~ Embedded system and application design
- Training in parallel programming (OpenMP, MPI, TBB, CUDA, OpenCL...)



Programming environments for big data processing on modern parallel architectures
 HPC Project Ronan KERYELL 9 / 101

Efficient big data architectures (I)

- Massive parallelism
- Use right processing elements for the right tasks (efficiency)
- Avoid moving data (expensive, slow)
- Avoid storing data: on-the-fly processing of data to be correlated
- Use memory hierarchies
- Distributed computing with smart routers




Programming environments for big data processing on modern parallel architectures
 HPC Project Ronan KERYELL 10 / 101

Hardware architectures

Outline

- 1 Hardware architectures
 - Classical multicores
 - GPU
 - MP-SoC underworlds
- 2 Software environments
 - Programming challenges
 - Multicores
 - GPU
 - Application libraries
- 3 Par4All
 - GPU code generation
 - Code generation for SCMP
- 4 Conclusion



Programming environments for big data processing on modern parallel architectures
 HPC Project Ronan KERYELL 11 / 101

Outline

- Hardware architectures
 - Classical multicores
 - GPU
 - MP-SoC underworlds
- Software environments
 - Programming challenges
 - Multicores
 - GPU
 - Application libraries
- Par4All
 - GPU code generation
 - Code generation for SCMP
- Conclusion

Programming environments for big data processing on modern parallel architectures
HPC Project Ronin KERVELL 12 / 101

Intel Nehalem (I)

- x86-64 new microarchitecture
- With on-chip memory controllers like AMD Opteron: ↗ bandwidth, ↘ latency
- SSE4.2 instruction set: strings (XML parser...), comparisons (data-mining), CRC & cryptography (protocols)
- Power consumption fine tuning with many sensors (power, temperature). Possible to speed up when less cores are used (turbo)
- Xeon X7560 (Beckton microarchitecture), 2010
 - 8 cores @ 2.27 GHz (+ turbo 2.666 GHz) + SMT 2 threads (HyperThreading) with 256 KB L2 cache/core, 32D+32I KB cache/core
 - 24 MB L3 cache
 - 4 Quick Path Interconnects (QPI) @ 6.4 GT/s to play Lego with processors & accelerators (≈HyperTransport d'AMD)

Programming environments for big data processing on modern parallel architectures
HPC Project Ronin KERVELL 13 / 101

Intel Nehalem (II)

- 4 DDR3-1333 MHz memory controllers
- 130 W
- \$3692 March 30, 2011
- 2.3 Gtr 45 nm 684 mm²

Programming environments for big data processing on modern parallel architectures
HPC Project Ronin KERVELL 14 / 101

Intel Nehalem (III)

Programming environments for big data processing on modern parallel architectures
HPC Project Ronin KERVELL 15 / 101

Intel Nehalem (IV)

Up to 8 sockets: Node Controllers optional
Larger than 8 sockets: Node Controllers required

EX = Intel Xeon processor 7500 series

Programming environments for big data processing on modern parallel architectures
HPC Project Ronin KERVELL 16 / 101

AMD Opteron 6180 (2011)

- x86-compatible instruction set
- 64-bit mode that double also register numbers
- Out-of-order superscalar execution with 9 instructions/cycle
 - 3 integer instructions
 - 3 address generators
 - 3 floating computation operators (+, *, memory)
- Opteron 6180 SE, 02/2011, \$1514
 - 12 cores @ 2.5 GHz, 512 KB/core L2 cache
 - 2 x 6 MB L3 cache
 - 21 GB/s DDR3 memory
 - 4 x 6.4 GT/s HyperTransport
 - 1.8 Gtr 45 nm 692 mm²
 - 140 W

Programming environments for big data processing on modern parallel architectures
HPC Project Ronin KERVELL 17 / 101

IBM Power 7 (2010)

- 4 chips per quad-chip module
- 8 cores per chip @ 4.0 GHz (4.25 GHz in TurboCore mode with 4 cores)
- 1.2 Gtr 45 nm SOI process, 567 mm²
- 4 SMT threads per core
- 321+32D kB L1 cache/core
- 256 kB L2 cache/core
- 32 MB L3 cache in eDRAM
- 100 GB/s DDR3 memory
- 12 execution units per core:
 - ▶ 2 fixed-point units
 - ▶ 2 load/store units
 - ▶ 4 double-precision floating-point units
 - ▶ 1 vector unit supporting VSX
 - ▶ 1 decimal floating-point unit
 - ▶ 1 branch unit
 - ▶ 1 condition register unit

Programming environments for big data processing on modern parallel architectures
HPC Project Ronan Keryell 18 / 101

Outline

- 1 Hardware architectures
 - Classical multicores
 - GPU
 - MP-SoC underwords
- 2 Software environments
 - Programming challenges
 - Multicores
 - GPU
 - Application libraries
- 3 Par4All
 - GPU code generation
 - Code generation for SCMP
- 4 Conclusion

Programming environments for big data processing on modern parallel architectures
HPC Project Ronan Keryell 19 / 101

Off-the-shelf AMD/ATI Radeon HD 6970 GPU

- 2.64 billion 40nm transistors
- 1536 stream processors @ 880 MHz, 2.7 TFLOPS SP, 675 GFLOPS DP
- + External 1 GB GDDR5 memory 5.5 Gt/s, 176 GB/s, 384b GDDR5
- 250 W on board (20 idle), PCI Express 2.1 x16 bus interface
- OpenGL, OpenCL
- ⊞ Radeon HD 6990 double chip card

More integration:

- Llano APU (FUSION Accelerated Processing Unit) : x86 multicore + GPU 32nm, OpenCL

Programming environments for big data processing on modern parallel architectures
HPC Project Ronan Keryell 20 / 101

Radeon HD 6870 — big picture

Programming environments for big data processing on modern parallel architectures
HPC Project Ronan Keryell 21 / 101

Off-the-shelf nVidia Tesla Fermi M2090 & GTX580

- GF110: 3 billion 40nm tr.
- 512 thread processors @ 1300 MHz, 1.3 TFLOPS SP, 666 GFLOPS DP
- + External 6 GB GDDR5 ECC memory 3.7 Gt/s, 177 GB/s. Less if using ECC

- 247 W on board PCI Express 2.1 x16 bus interface
- OpenGL, OpenCL, CUDA

Programming environments for big data processing on modern parallel architectures
HPC Project Ronan Keryell 22 / 101

GF100 Stream Multiprocessor

Programming environments for big data processing on modern parallel architectures
HPC Project Ronan Keryell 23 / 101

Outline

- 1 Hardware architectures
 - Classical multicores
 - GPU
 - MP-SoC underworlds
- 2 Software environments
 - Programming challenges
 - Multicores
 - GPU
 - Application libraries
- 3 Pa4All
 - GPU code generation
 - Code generation for SCMP
- 4 Conclusion

Programming environments for big data processing on modern parallel architectures
HPC Project Ronin Kerivell 24 / 101

ARM yourself

- Do some computations where the caplors are...
- Smartphone and other sensor networks
- Trade-off between communication energy and inside/remote computations
- Texas Instrument OMAP4470 announced on 2011/06/02
 - ▶ 2 ARM Cortex-A9 MPCores @ 1.8GHz with Neon vector instructions
 - ▶ 2 ARM Cortex-M3 cores (low-power and real-time responsiveness, multimedia, avoiding to wake up the Cortex-A9...)
 - ▶ **SGX544 graphics core with OpenCL 1.1 support**, with 4 USSE2 core @ 384 MHz producing each 4 FMAD/cycle: 12.3 GFLOPS
 - ▶ 2D graphics accelerator
 - ▶ 3 HD displays and up to QXGA (2048x1536) resolution + stereoscopic 3D
 - ▶ Dual-channel, 466 MHz LPDDR2 memory

Programming environments for big data processing on modern parallel architectures
HPC Project Ronin Kerivell 25 / 101

ARM yourself

- ~ Current course to have non-x86 servers based on ARM...
- ∃ Experiments on low power clusters
- Think to evaluate power consumption on your application

Programming environments for big data processing on modern parallel architectures
HPC Project Ronin Kerivell 26 / 101

Tilera TilePro64

Programming environments for big data processing on modern parallel architectures
HPC Project Ronin Kerivell 27 / 101

Tilera TilePro64

- Processing power with good network interfaces
 - ▶ Interactive data analysis
 - ▶ Video/audio codec
 - ▶ DPI, IDS, IPsec...
- 8x8 processors with SMP, partitionable
- 32-bit VLIW with SIMD mode
- 700-866 MHz: 443 GOPS 8 bits, 23 W
- 64 DDR2 controllers 25.6 GB/s
- 2x 10 GbE XAUI + IP session hash distribution
- SMP Linux or bare bone per tile
- C/C++ gcc compiler
- TMC library for hardcore support (2D network...)
- Eclipse support with graphical simulator
- OpenMP

Programming environments for big data processing on modern parallel architectures
HPC Project Ronin Kerivell 28 / 101

MPPA de Kalray (the French touch!)

- processors)
 - Shared memory and NoC with DMA
 - 28 nm CMOS technology, ≈ 5 W @ 400 MHz
 - FPU 32/64 bits IEEE 754: 205 GFLOPS SP, 1 TOPS 16 bits
 - 2x 64-bit DDR3 memory controllers for high bandwidth main memory transfers
 - 2x 40 Gb/s or 8x 10 Gb/s Ethernet controller
- 256 VLIW processors per chip 16 clusters of 16

Programming environments for big data processing on modern parallel architectures
HPC Project Ronin Kerivell 29 / 101

Hardware architectures MP-SoC underworks

MPPA de Kalray (the French touch!) (II)

- 2× 8-lane PCI Express Gen 3
- 4× 4–8-lane Interlaken interfaces for multi-MPPA chip system integration (8 MPPA/PCIe board) or connection to external FPGAs, I/O...
- Linux or bare metal with AccessCore library
- Multi-core compiler (gcc 4.5), simulator, debugger (gdb), profiler
- Eclipse IDE
- Programming from high-level C-based language
- AccessCore library

Programming environments for big data processing on modern parallel architectures
HPC Project Ronan Keryell 30 / 101

Hardware architectures MP-SoC underworks

FPGA (I)

- Field-programmable gate array with bitstream configuration in memory
- Xilinx Virtex 7
 - ▶ 2M logic cells (6-LUT), 28 nm
 - ▶ 85Mb block RAM
 - ▶ 5280 DSP slices (6.7 TFMA/s)
 - ▶ 96 transceiver @ 28Gb/s: 2.8 Tb/s
 - ▶ PCIe gen3 ×8
 - ▶ 1200 pins
- Radar, communications, HPC, datamining, bioinformatics...

Programming environments for big data processing on modern parallel architectures
HPC Project Ronan Keryell 31 / 101

Hardware architectures MP-SoC underworks

FPGA (II)

- VHDL-to-bitstream compiler... but **hard work**
- Dynamic partial reconfiguration
- C-to-VHDL compilers
 - ▶ Riverside Optimizing Compiler for Configurable Computing (ROCCC): open source
 - ▶ Impulse C
 - ▶ Catapult C
 - ▶ Cynthesizer
 - ▶ ...

Programming environments for big data processing on modern parallel architectures
HPC Project Ronan Keryell 32 / 101

Hardware architectures MP-SoC underworks

Convey HC-1^{ex}

- Intel Xeon quad-core @2.13 GHz with 128 GB
- 4 Xilinx Virtex 6 LX760 FPGA with 128 GB DDR2
- 1520 W in 3U rack
- Linux
- Various instruction sets ("personality")
- Fortran/C/C++ vectorizing compiler replacing complex instructions by FPGA vectorized implementation
- Possible to design its own instruction set (ROCCC...)
- 401× speed-up on Smith-Waterman algorithm compared to x86

Programming environments for big data processing on modern parallel architectures
HPC Project Ronan Keryell 33 / 101

Hardware architectures MP-SoC underworks

Anton computer from D.E. Shaw

- 1 Create a hedge fund
- 2 Earn a lot of money
- 3 Spend it by creating a 500+ people start-up in bioinformatics
- 4 Build from scratch (ASIC) a computer for solving special issues

Programming environments for big data processing on modern parallel architectures
HPC Project Ronan Keryell 34 / 101

Hardware architectures MP-SoC underworks

PacketShader


- Use Linux PC + GPU as a router with processing power
 - ▶ 2 4-core Nehalem @2.66 GHz + 4 10 GigE NIC + 2 GPU GTX480
- 40 Gb/s routing even with 64-byte packets
- 8–20 Gb/s IPsec tunneling
- Linux IP stack to slow → own raw device driver for Intel 82598/82599 NIC
 - ▶ Extract Ethernet flow into user mode
 - ▶ Huge recycled packet buffers
 - ▶ Batch processing to amortize overhead
 - ▶ NUMA-aware processing: packets processed by NIC-local GPU in its RAM, aligned in cache
- Current bottleneck: IOH chipset

<http://shader.kaist.edu/packetshader>

Programming environments for big data processing on modern parallel architectures
HPC Project Ronan Keryell 35 / 101

OpenFlow


- Need for advanced routing with computing capabilities
- PC or Tilera-like with few 10+ Gb/s Ethernet or Infiniband: bandwidth maybe OK but not enough links
- Open standard to interact with existing routers & switch: OpenFlow <http://www.openflow.org>
- Possible to extract *minimal* flows to feed PC/GPU/MP-SoC/FPGA accelerators and reinject results in the router
- OpenFlow implemented by many router companies
- Possible with non-OpenFlow but less flexible and non portable



Programming environments for big data processing on modern parallel architectures
HPC Project Ronin KERVELL 38 / 101

Outline


- 1 Hardware architectures
 - Classical multicores
 - GPU
 - MP-SoC underwords
- 2 Software environments
 - Programming challenges
 - Multicores
 - GPU
 - Application libraries
- 3 Par4All
 - GPU code generation
 - Code generation for SCMP
- 4 Conclusion



Programming environments for big data processing on modern parallel architectures
HPC Project Ronin KERVELL 37 / 101

Outline

- 1 Hardware architectures
 - Classical multicores
 - GPU
 - MP-SoC underwords
- 2 Software environments
 - Programming challenges
 - Multicores
 - GPU
 - Application libraries
- 3 Par4All
 - GPU code generation
 - Code generation for SCMP
- 4 Conclusion




Programming environments for big data processing on modern parallel architectures
HPC Project Ronin KERVELL 38 / 101

Parallel application dwarfs

- <http://view.eecs.berkeley.edu/> : « The Landscape of Parallel Computing Research: A View From Berkeley »
- Try to capture typical examples to analyze and design new architectures & applications


	Dwarf	Performance Limit: Memory Bandwidth, Memory Latency, or Computation?
1	Dense Matrix	Computationally limited
2	Sparse Matrix	Currently 50% computation, 50% memory BW
3	Spectral (FFT)	Memory latency limited
4	N-Body	Computationally limited
5	Structured Grid	Currently more memory bandwidth limited
6	Unstructured Grid	Memory latency limited
7	MapReduce	Problem dependent
8	Combinational Logic	CRC problems BW, crypto problems computationally limited
9	Graph traversal	Memory latency limited
10	Dynamic Programming	Memory latency limited
11	Backtrack and Branch-Bound	?
12	Construct Graphical Models	?
13	Finite State Machine	Nothing helps!



Programming environments for big data processing on modern parallel architectures
HPC Project Ronin KERVELL 39 / 101

Extracting parallelism in applications...

- The implicit attitude
 - ▶ Hardware: massively superscalars processors
 - ▶ Software: auto-parallelizing compilers
- The (±) explicit attitude
 - ▶ Languages (± extensions): OpenMP, UPC, HPF, Co-array Fortran (F-), Fortran 2008, X10, Chapel, Fortress, Matlab, SciLab, Octave, Mapple, LabView, nVidia CUDA, AMD/ATI Stream (Brook+, Cal), OpenCL, HMPP, *insert your own preferred language here*...
 - ▶ Framework: MapReduce, Hadoop...
 - ▶ Libraries: application-oriented (mathematics, coupling...), parallelism (MPI, concurrency *threads*, SPE/MFC on Cell...), Multicore Association MCAPI, objects (parallel STL, TBB, Ct...)




Programming environments for big data processing on modern parallel architectures
HPC Project Ronin KERVELL 40 / 101

... but multidimensional heterogeneity!

Welcome into Parallel Hard-Core Real Life 2.0!

- Heterogeneous execution models
 - ▶ Multicore SMP ± coupled by caches
 - ▶ SIMD instructions in processors (Mnem, VMX, SSE4.2, 3DNow!, LRBn...)
 - ▶ Hardware accelerators (MIMD, MISD, SIMD, SIMT, FPGA...)
- New heterogeneous memory hierarchies
 - ▶ Classic caches/physical memory/disks
 - ▶ Flash SSD is a new-comer to play with
 - ▶ NUMA (*Non Uniform Memory Access*) : sockets-attached memory banks, remote nodes...
 - ▶ Peripherals attached to sockets : NUPA (*Non Uniform Peripheral Access*). GPU on PCIe x16 in this case...
 - ▶ If non-shared memory: remote memory, remote disks...
 - ▶ Inside GPU : registers, local memory, shared memory, constant memory, texture cache, processor grouping, locked physical pages, host memory access...
- Heterogeneous communications
 - ▶ Anisotropic networks
 - ▶ Various protocols

⚠ Several dimensions to cope with at the same time



Programming environments for big data processing on modern parallel architectures
HPC Project Ronin KERVELL 41 / 101

Outline

- 1 Hardware architectures
 - Classical multicores
 - GPU
 - MP-SoC underworlds
- 2 Software environments
 - Programming challenges
 - Multicores
 - GPU
 - Application libraries
- 3 Par4All
 - GPU code generation
 - Code generation for SCMP
- 4 Conclusion

Programming environments for big data processing on modern parallel architectures
HPC Project Ronan Keryell 42 / 101

OpenMP

- « Le multithread pour les nuls » ☺
- Vise machines à mémoire partagée
- Sauf si programmation système compliquée, pas besoin de faire de la programmation de threads explicites
- Idée : saupoudrer un programme de directives pour aider compilateur à paralléliser
- Philosophie : #pragmatisme avec une certaine élégance esthétique
- ⚠ Si pas de directives, pas de parallélisme exploité (*a priori*)
- ⚠⚠⚠ Directive ≡ déclaration sur l'honneur
- Langages supportés : Fortran et C/C++

<http://openmp.org>

Programming environments for big data processing on modern parallel architectures
HPC Project Ronan Keryell 43 / 101

Modèle d'exécution d'OpenMP (I)

<http://openmp.org/wp/openmp-specifications>

- Exécution parallèle SPMD basée sur le *fork/join*

- Création de thread implicite ou explicite avec des *directives*

Programming environments for big data processing on modern parallel architectures
HPC Project Ronan Keryell 44 / 101

Modèle d'exécution d'OpenMP (II)

- Astuce : un programme OpenMP peut être exécuté comme
 - ▶ Programme séquentiel
 - ▶ Programme parallèle
- ↪ portabilité, coût de sortie nul ! ☺
- Threads créées dans des sections `parallel` et stoppées à la fin avec une barrière
- Constructions de synchronisation et dans bibliothèque
- Contrôle de l'environnement d'exécution par variables d'environnement et fonctions de bibliothèque

Programming environments for big data processing on modern parallel architectures
HPC Project Ronan Keryell 45 / 101

Exemple

```
#pragma omp parallel default(none) \
shared(n,x,y) private(i)
{
  /* Ceci s'exécute sur plusieurs threads en // */
  #pragma omp for
  for (i=0; i<n; i++)
    // Les itérations sont réparties sur les threads
    x[i] += y[i];
  // Synchronisation implicite ici
} /* End of parallel region */
// Synchronisation implicite ici
```

Programming environments for big data processing on modern parallel architectures
HPC Project Ronan Keryell 46 / 101

Task en OpenMP 3.0

- Rajout de la notion de tâches explicites ≡ bout de programme exécuté sur une thread
- Tâche créée dans thread par construction `task` (TBB, Cilk...)

```
#pragma omp parallel
{
  #pragma omp single private(p)
  {
    p = listhead;
    while (p) {
      #pragma omp task
      {
        process (p);
      }
      p=next(p);
    }
  }
}
```

- Extensions en vue (target...)

Programming environments for big data processing on modern parallel architectures
HPC Project Ronan Keryell 47 / 101

Multimedia SIMD extensions (I)

- Multimedia and telecommunication market ↗ ↘
 - GIF image: 8-bit pixels
 - RGB_a true-color image: 4 × 8 bit pixels
 - Phone A- or μ-law : 8-bit samples
 - CD-quality sound : 2 × 16-bit samples
- In general purpose processors, under-used transistors on these applications (double precision multiplier...) ☹
- Idea : 128-bit data viewed as independent vectors of 16 independent 8-bit elements or 8 16-bit or... 2 DP-float or 4 SP-float
- ~ Add SIMD to general purpose processors (i860, SSE 4.2 Core i7, AMD 3Dnow!, VMX Power, ARM, SPARC...) and Cell for numerical computing, strings, data compaction, cryptography...
- SSE4.2 at 3.2 GHz : 2 128-bit operations/cycle ~ 819 GOP1b/s, 102 GOP8b/s *per core!* ☺

Programming environments for big data processing on modern parallel architectures
HPC Project Ronin KERVELL 48 / 101

Multimedia SIMD extensions (II)

- Programming
 - Quite complex ☹ because many *ad hoc* instructions, saturation arithmetic...
 - Just do it... in assembly language ☹
 - Intrinsics functions in C/C++ (GCC, Intel...)
 - C/C++ extensions with new vector data types (GCC, Intel...)
 - Auto-vectorizing (IBM x1c, Intel, GCC, generic tools such as PIPS)
 - Use already optimized libraries
- Available on all machines ☺

Programming environments for big data processing on modern parallel architectures
HPC Project Ronin KERVELL 49 / 101

Message Passing Interface (MPI) (I)

- Le passage de message pour les nuls ! ☹
- Bibliothèque de fonctions de communication disponible pour de nombreux langages et systèmes d'exploitation
- Portabilité et nivellement par le bas : programmation de SMP aussi en MPI...
- Ressources
 - http://en.wikipedia.org/wiki/Message_Passing_Interface
 - MPI Forum <http://www.mpi-forum.org>
 - MPICH : A Portable Implementation of MPI <http://www-unix.mcs.anl.gov/mpi/mpich/>
 - LAM / MPI Parallel Computing <http://www.mpi.nd.edu/lam/>
 - MPE Graphics-Scalable X11 Graphics in MPI <http://www-fp.mcs.anl.gov/~lusk/papers/mpe/>
 - Livre « Using MPI. Portable Parallel Programming with the Message-Passing Interface » <http://www-unix.mcs.anl.gov/mpi/usingmpi/>

Programming environments for big data processing on modern parallel architectures
HPC Project Ronin KERVELL 50 / 101

Message Passing Interface (MPI) (II)

- Faute de temps je n'utilise plus mes transparents trop complets <http://enstb.org/~keryell/cours/MR2/IAHP/MPI>
- J'utilise « An Introduction to MPI Parallel Programming with the Message Passing Interface » de William Gropp & Ewing Lusk <http://www-unix.mcs.anl.gov/mpi/tutorial/mpiintro/MPIIntro.PPT>

Programming environments for big data processing on modern parallel architectures
HPC Project Ronin KERVELL 51 / 101

Thread Building Blocks (TBB) (I)

<http://en.wikipedia.org/wiki/TBB> from Intel

- Template library (*à la* STL)
- Open and commercial Versions
- Algorithms (for, reduce, pipeline, scan...), containers, memory allocators, mutual exclusion, atomic operations, schedulers, profiling...
- Work stealing between tasks
- Orthogonal to OpenMP & MPI

Programming environments for big data processing on modern parallel architectures
HPC Project Ronin KERVELL 52 / 101

Thread Building Blocks (TBB) (II)

```

1 class ApplyFoo {
2     float* const my_a;
3 public:
4     ApplyFoo(float* a) : my_a(a) {};
5     void operator() (const tbb::blocked_range<size_t>&r) const {
6         for (size_t i=r.begin(); i != r.end(); ++i)
7             Foo (my_a[i]);
8     }
9 }
10
11 void ParallelApplyFoo(float a[], size_t n) {
12     tbb::parallel_for (
13         tbb::blocked_range<size_t>(0,n),
14         ApplyFoo(a),
15         tbb::auto_partitioner ()
16     );
17 }

```

- Need a deep code restructuring if an application is not in a STL spirit

Programming environments for big data processing on modern parallel architectures
HPC Project Ronin KERVELL 53 / 101

Outline

- 1 Hardware architectures
 - Classical multicores
 - GPU
 - MP-SoC underworlds
- 2 Software environments
 - Programming challenges
 - Multicores
 - GPU
 - Application libraries
- 3 Par4All
 - GPU code generation
 - Code generation for SCMP
- 4 Conclusion

Programming environments for big data processing on modern parallel architectures
HPC Project
Ronan Keryell
54 / 101

Basic GPU programming model

A sequential program on a host launches computational-intensive kernels on a GPU

- Allocate storage on the GPU
- Copy-in data from the host to the GPU
- Launch the kernel on the GPU
- The host waits...
- Copy-out the results from the GPU to the host
- Deallocate the storage on the GPU

Programming environments for big data processing on modern parallel architectures
HPC Project
Ronan Keryell
55 / 101

GPU execution model

Programming environments for big data processing on modern parallel architectures
HPC Project
Ronan Keryell
56 / 101

From hardware constraints to programming style (I)

- GPU computes fast but connected to CPU with slow PCI link ~>
 - ▶ Avoid exchanging too much data between CPU & GPU or... compute on the CPU ☺
 - ▶ Possible to overlap communications with computations (more complex programming)
- Many SIMD engines (multiprocessors) ~> at least as much blocks of threads
- Memory hierarchy is quite complex and... visible!
 - ▶ Use (quite limited ☹) local registers by recycling local data
 - ▶ Memory is accessed in huge lines ~> program to use all the elements of the line
 - ▶ If not possible, try to reorganize data in the shared memory around read/write (matrix transposition...)
 - ▶ Recently added caches help too

Programming environments for big data processing on modern parallel architectures
HPC Project
Ronan Keryell
57 / 101

From hardware constraints to programming style (II)

- ▶ Memory is far far away (800+ cycles) ~> use a lots of thread per block (but limited resources reduce block numbers) to overlap memory access with other computations
- ▶ Computing is fast, memory is slow. Rethink algorithms...
- SIMD machine, only one control flow ~> predicated

```

1 if (cond[i])
  b[i] = a[i];
3 else
  b[i] = -a[i] + 1;
    
```

- ▶ Some hardware optimizations if in a SIMD warp there is no execution ~> if possible sort false/true elements

Programming environments for big data processing on modern parallel architectures
HPC Project
Ronan Keryell
58 / 101

Programmation CUDA (I)

- Data-parallel extension to a C++ subset
- Target nVidia GPU and x86 multicores
- 2-level parallelism: threads in blocks of threads + block-tiling
- In a block of threads : communication through shared memory and synchronization via `__syncthreads()`
- Complex heterogeneous memory layout (GPU...)

```

__global__ void
add_matrix_gpu(float *a, float *b, float *c, int N) {
  int i=blockIdx.x*blockDim.x+threadIdx.x;
  int j=blockIdx.y*blockDim.y+threadIdx.y;
  int index =i+j*N;

  if ( i < N && j < N )
    c[index]=a[index]+b[index];
}

void main() {
  float ha[N][N], hb[N][N], hc[N][N];
    
```

Programming environments for big data processing on modern parallel architectures
HPC Project
Ronan Keryell
59 / 101

Programmation CUDA (II)

```

/* Allocate array on the GPU with cudaMalloc */
float *a, *b, *c;
cudaMalloc((void **) &a, sizeof(float)*N*N);
cudaMalloc((void **) &b, sizeof(float)*N*N);
cudaMalloc((void **) &c, sizeof(float)*N*N);

cudaMemcpy(a, ha, sizeof(float)*N*N, cudaMemcpyHostToDevice);
cudaMemcpy(b, hb, sizeof(float)*N*N, cudaMemcpyHostToDevice);

// Describe iteration tiling (2D strip-mining)
dim3 dimBlock (blocksize,blocksize);
dim3 dimGrid (N/dimBlock.x,N/dimBlock.y);
add_matrix_gpu<<<dimGrid,dimBlock>>>(a,b,c,N);
cudaMemcpy(c, hc, sizeof(float)*N*N, cudaMemcpyDeviceToHost);
}

```

- Need some heavy code restructuring
- ∃ other version: CUDA driver, similar to OpenCL

OpenCL (I)

- Language based on a C99 subset
- Started by Apple to *unify* parallel use (multicores, GPGPU...)
 - ↳ similar to OpenGL & OpenGL
- Followed by AMD/ATI and nVidia
- Data-parallelism and control-parallelism (1–3-dimensions) according to targets
- Kernel oriented computations on streams
- Complex split memory model (GPGPU...) but CPU compliant too
- New types (vectors, images...)

OpenCL (II)

```

/* This kernel computes FFT of length 1024.
The 1024 length FFT is decomposed into calls to a radix 16
function, another radix 16 function and then a radix 4 function */
__kernel void fft1D_1024 (__global float2 *in, __global float2 *out,
                        __local float *sMemx, __local float *sMemy) {
    int tid = get_local_id(0);
    int blockIdx = get_group_id(0) * 1024 + tid;
    float2 data[16];
    // starting index of data to/from global memory
    in = in + blockIdx; out = out + blockIdx;
    globalLoads(data, in, 64); // coalesced global reads
    fftRadix16Pass(data); // in-place radix-16 pass
    twiddleFactorMul(data, tid, 1024, 0); // in-place twiddle factor multiplication
    localShuffle(data, sMemx, sMemy, tid,
                ((tid & 15) * 65) + (tid >> 4));
    fftRadix16Pass(data); // in-place radix-16 pass
    twiddleFactorMul(data, tid, 64, 4); // twiddle factor multiplication
    localShuffle(data, sMemx, sMemy, tid,
                ((tid >> 4) * 64) + (tid & 15));
    // four radix-4 function calls
    fftRadix4Pass(data); fftRadix4Pass(data + 4);
    fftRadix4Pass(data + 8); fftRadix4Pass(data + 12);
}

```

OpenCL (III)

```

// coalesced global writes
globalStores(data, out, 64);
}
[...]
// create a compute context with CPU device
context = clCreateContextFromType(CL_DEVICE_TYPE_GPU);
// create a work-queue
queue = clCreateWorkQueue(context, NULL, NULL, 0);
// allocate the buffer memory objects
memobjs[0] = clCreateBuffer(context, CL_MEM_READ_ONLY |
                           CL_MEM_COPY_HOST_PTR,
                           sizeof(float)*2*num_entries, srcA);
memobjs[1] = clCreateBuffer(context, CL_MEM_READ_WRITE,
                           sizeof(float)*2*num_entries, NULL);
// create the compute program
program = clCreateProgramFromSource(context, 1,
                                   &fft1D_1024_kernel_src, NULL);
// build the compute program executable
clBuildProgramExecutable(program, CL_MEM_READ_WRITE,
                        clBuildProgramExecutable(program, false, NULL, NULL));
// create the compute kernel
kernel = clCreateKernel(program, "fft1D_1024");
// create N-D range object with work-item dimensions
global_work_size[0] = n;
local_work_size[0] = 64;

```

OpenCL (IV)

```

range = clCreateNDRangeContainer(context, 0, 1,
                                global_work_size, local_work_size);
// set the args values
clSetKernelArg(kernel, 0, (void *)&memobjs[0], sizeof(cl_mem), NULL);
clSetKernelArg(kernel, 1, (void *)&memobjs[1], sizeof(cl_mem), NULL);
clSetKernelArg(kernel, 2, NULL,
               sizeof(float)*(local_work_size[0]+1)*16, NULL);
clSetKernelArg(kernel, 3, NULL,
               sizeof(float)*(local_work_size[0]+1)*16, NULL);
// execute kernel
clExecuteKernel(queue, kernel, NULL, range, NULL, 0, NULL);

```

- Need a lot of code restructuring

CUDA or OpenCL?

CUDA

- Appeared first
- Language basis not well defined: C++ like, rather C89 and not C99
- Painful to translate C99 to C89 and keeping clean sources
- Rather limited to 2D threads
- nVidia GPU only

OpenCL

- Standard backed by many companies
- C99 based ↳ clean
- 3D threads with less constraints
- More verbose API (kernel call...)
- Kernel source code outside of host source: more complex
- Fast spreading in embedded computing world (MP-SoC)

Outline

- 1 Hardware architectures
 - Classical multicores
 - GPU
 - MP-SoC underworlds
- 2 Software environments
 - Programming challenges
 - Multicores
 - GPU
 - Application libraries
- 3 Par4All
 - GPU code generation
 - Code generation for SCMP
- 4 Conclusion

Programming environments for big data processing on modern parallel architectures
HPC Project
Ronan KERYELL
68 / 101

Bibliothèques mathématiques (I)

- Existent pour différentes architectures !
- Souvent bibliothèques constructeurs optimisées pour leur machines
 - ▶ Intel Math Kernel Library (MKL)
 - ▶ AMD Performance Library (↪ Framework libre)
- FFT : FFTW (*Fastest Fourier Transform in the West*, en C généré par du OCaml)...
- Beaucoup d'algèbre linéaire
 - ▶ BLAS (*Basic Linear Algebra Subprograms*) et PBLAS
 - ▶ LAPACK (*Linear Algebra PACKage*)
 - ▶ ScaLAPACK : version SPMD avec MPI
 - ▶ SuperLU : solution directe de gros systèmes creux
 - ▶ PETSc (*Portable, Extensible Toolkit for Scientific Computation*) : large spectre, au dessus de MPI

Programming environments for big data processing on modern parallel architectures
HPC Project
Ronan KERYELL
67 / 101

Outline

- 1 Hardware architectures
 - Classical multicores
 - GPU
 - MP-SoC underworlds
- 2 Software environments
 - Programming challenges
 - Multicores
 - GPU
 - Application libraries
- 3 Par4All
 - GPU code generation
 - Code generation for SCMP
- 4 Conclusion

Programming environments for big data processing on modern parallel architectures
HPC Project
Ronan KERYELL
68 / 101

Use the Source, Luke...

Hardware is moving quite (too) fast but...

What has survived for 50+ years?
Fortran programs...

What has survived for 40+ years?
IDL, Matlab, Scilab...

What has survived for 30+ years?
C programs, Unix...

- A lot of legacy code could be pushed onto parallel hardware (accelerators) with automatic tools...
- Need automatic tools for source-to-source transformation to leverage existing software tools for a given hardware
- Not as efficient as hand-tuned programs, but quick production phase

Programming environments for big data processing on modern parallel architectures
HPC Project
Ronan KERYELL
69 / 101

Not reinventing the wheel... No NIH syndrome please!

Want to create your own tool?

- House-keeping and infrastructure in a compiler is a **huge** task
- Unreasonable to begin yet another new compiler project...
- Many academic Open Source projects are available...
- ...But customers need products ☹
- ↪ Integrate your ideas and developments in existing project
- ...or buy one if you can afford (ST with PGI...) ☹
- Some projects to consider
 - ▶ Old projects: gcc, PIPS... and many dead ones (SUIF...)
 - ▶ But new ones appear too: LLVM, RoseCompiler, Cetus...

Par4All

- ↪ Funding an initiative to industrialize Open Source tools
- Use source-to-source tool to be more target-independent
- PIPS is the first project to enter the Par4All initiative

Programming environments for big data processing on modern parallel architectures
HPC Project
Ronan KERYELL
70 / 101

Current PIPS usage


Developed for 23 years (...) @ Mines ParisTech & Télécom Bretagne, mainly

- Automatic parallelization (Par4All C & Fortran to OpenMP)
- Distributed memory computing with OpenMP-to-MPI translation [STEP project]
- Generic vectorization for SIMD instructions (SSE, VMX, Neon, CUDA, OpenCL...) (SAC project) [SCALOPEs]
- Parallelization for embedded systems [SCALOPEs]
- Compilation for hardware accelerators (Ter@PIX, SPoC, SIMD, FPGA, MPPA, P2012...) [FREIA, SCALOPEs, SMECY]
- High-level hardware accelerators synthesis generation for FPGA [PHRASE, CoMap]
- Reverse engineering & decompiler (reconstruction from binary to C)
- Genetic algorithm-based optimization [Luxembourg university+TB]
- Code instrumentation for performance measures
- GPU with CUDA & OpenCL [TransMedi@, FREIA, OpenGPU]

Programming environments for big data processing on modern parallel architectures
HPC Project
Ronan KERYELL
71 / 101

Par4All usage

- Generate from sequential C, Fortran & Scilab code
 - ▶ OpenMP for SMP
 - ▶ CUDA for nVidia GPU
 - ▶ OpenCL for GPU & ST Platform 2012 (on-going)
 - ▶ Code for various accelerators [SMECY], Kalray [SIMILAN]... (on-going)
 - ▶ SCMP task programs for SCMP machine from CEA and for... cloud computing ☺




Programming environments for big data processing on modern parallel architectures
HPC Project Ronin KERVELL 72 / 101

Par4All ≡ PyPS scripting in the backstage (I)

- PIPS is a great tool-box to do source-to-source compilation
- ...but not really usable by λ end-user ☹
- \rightsquigarrow Development of Par4All
- Add a user compiler-like infrastructure

\rightsquigarrow p4a script as simple as

 - ▶ `p4a --openmp toto.c -o toto`
 - ▶ `p4a --cuda toto.c -o toto -lm`
- Be multi-target
- Apply some adaptative transformations
- Up to now PIPS was scripted with a special shell-like language: `trips`
- Not enough powerful (not a programming language)
- Develop a SWIG Python interface to PIPS phases and interface



Programming environments for big data processing on modern parallel architectures
HPC Project Ronin KERVELL 73 / 101

Par4All ≡ PyPS scripting in the backstage (II)


- ▶ All the power of a widely spread real language
- ▶ Automate with introspection through the compilation flow
- ▶ Easy to add any glue, pre-/post-processing to generate target code

Overview

```

graph LR
    subgraph Sequential_source_code [Sequential source code]
        f[f]
        c[c]
    end
    subgraph Par4All
        subgraph Preprocessor
            f --> P[Preprocessor]
        end
        subgraph PyPS_PIPS [PyPS PIPS]
            P --> PIPS[PyPS PIPS]
        end
        subgraph Postprocessor
            PIPS --> Post[Postprocessor]
        end
        subgraph P4A_Accel_runtime [P4A Accel runtime]
            Post --> P4A[P4A Accel runtime]
        end
        subgraph Back_end_compilers [Back-end compilers]
            P4A --> gcc[gcc, icc, ...]
            P4A --> nvcc[nvcc]
        end
    end
    subgraph Parallel_source_code [Parallel source code]
        gcc --> p4a_c[p4a.c]
        nvcc --> p4a_cu[p4a.cu]
    end
    subgraph Parallel_executables [Parallel executables]
        p4a_c --> OpenMP[OpenMP executu]
        p4a_cu --> CUDA[CUDA executet]
    end
  
```


- Invoke PIPS transformations
 - ▶ With different recipes according to generated stuff
 - ▶ Special treatments on kernels...
- Compilation and linking infrastructure: can use `gcc`, `icc`, `nvcc`, `nvcc+gcc`, `nvcc+icc`



Programming environments for big data processing on modern parallel architectures
HPC Project Ronin KERVELL 74 / 101

Par4All ≡ PyPS scripting in the backstage (III)


- House keeping code
- Fundamental: coloring and filtering some PIPS output, running cursor... ☹



Programming environments for big data processing on modern parallel architectures
HPC Project Ronin KERVELL 75 / 101

Coding rules


- Automatic parallelization is not magic
- Use abstract interpretation to « understand » programs
- Undecidable in the generic case (\approx halting problem)
- Quite easier for well written programs
- Develop a coding rule manual to help parallelization and... sequential quality!
 - ▶ Avoid useless pointers
 - ▶ Take advantage of C99 (arrays of non static size...)
 - ▶ Use higher-level C, do not linearize arrays...
 - ▶ Organize execution in cleaner loops expressing better parallelism
 - ▶ ...
- Prototype of coding rules report on-line on par4all.org



Programming environments for big data processing on modern parallel architectures
HPC Project Ronin KERVELL 76 / 101

Parallelization to OpenMP (I)

- The easy way... Already in PIPS
- Used to bootstrap the start-up with stage-0 investors ☹
- Indeed, we used only `bash`-generated `trips` at this time (2008, no PyPS yet), but needed a lot of bug squashing on C support in PIPS...



Programming environments for big data processing on modern parallel architectures
HPC Project Ronin KERVELL 77 / 101

OpenMP output sample

```

!$omp parallel do private(I, K, X)
C multiply the two square matrices of ones
DO J = 1, N
0016
!$omp parallel do private(K, X)
DO I = 1, N
0017
X = 0
0018
!$omp parallel do reduction(+-X)
DO K = 1, N
0019
X = X+A(I,K)*B(K,J)
0020
ENDDO
!$omp end parallel do
C(I,J) = X
0022
ENDDO
!$omp end parallel do
ENDDO
!$omp end parallel do

```

Programming environments for big data processing on modern parallel architectures
HPC Project Ronan Keryell 78 / 101

Outline

- Hardware architectures
 - Classical multicores
 - GPU
 - MP-SoC underworlds
- Software environments
 - Programming challenges
 - Multicores
 - GPU
 - Application libraries
- Par4All
 - GPU code generation
 - Code generation for SCMP
- Conclusion

Programming environments for big data processing on modern parallel architectures
HPC Project Ronan Keryell 79 / 101

Basic GPU execution model

A sequential program on a host launches computational-intensive kernels on a GPU

- Allocate storage on the GPU
- Copy-in data from the host to the GPU
- Launch the kernel on the GPU
- The host waits...
- Copy-out the results from the GPU to the host
- Deallocate the storage on the GPU

Generic scheme for other heterogeneous accelerators too

Programming environments for big data processing on modern parallel architectures
HPC Project Ronan Keryell 80 / 101

Scilab & Matlab

- Languages used for simulations, data analytics, pricing...
- Scilab/Matlab input : *sequential* or array syntax
- Compilation to C code
 - Side effect of MediaGPU ANR project...
 - Our COLD compiler is *not* Open Source
 - There is such Open Source compiler from hArtes European project written in... Scilab ☺
- Parallelization of the generated C code
- Use parallel runtime too
- Type inference to guess (crazy ☹) semantics
 - Heuristic: first encountered type is forever
- May get speedup $\gg 1000$ ☺
- Wild Cruncher product from HPC Project: x86+GPU appliance with nice interface
 - Scilab — mathematical model & simulation
 - Par4All — automatic parallelization
 - //Geometry — polynomial-based 3D rendering & modelling

Programming environments for big data processing on modern parallel architectures
HPC Project Ronan Keryell 81 / 101

Hyantes (I)

- Geographical application for data integration: library to compute neighbourhood population potential with scale control
- WildNode with 2 Intel Xeon X5670 @ 2.93GHz (12 cores) and a nVidia Tesla C2050 (Fermi), Linux/Ubuntu 10.04, gcc 4.4.3, CUDA 3.1
 - Sequential execution time on CPU: 30.355s
 - OpenMP parallel execution time on CPUs: 3.859s, speed-up: 7.87
 - CUDA parallel execution time on GPU: 0.441s, speed-up: 68.8
- With single precision on a HP EliteBook 8730w laptop (with an Intel Core2 Extreme Q9300 @ 2.53GHz (4 cores) and a nVidia GPU Quadro FX 3700M (16 multiprocessors, 128 cores, architecture 1.1)) with Linux/Debian/sid, gcc 4.4.5, CUDA 3.1:
 - Sequential execution time on CPU: 34.7s
 - OpenMP parallel execution time on CPUs: 13.7s, speed-up: 2.53
 - OpenMP emulation of GPU on CPUs: 9.7s, speed-up: 3.6
 - CUDA parallel execution time on GPU: 1.57s, speed-up: 24.2

Programming environments for big data processing on modern parallel architectures
HPC Project Ronan Keryell 82 / 101

Hyantes (II)

Original main C kernel:

```

void run(data_t xmin, data_t ymin, data_t xmax, data_t ymax, data_t step, d
town pt [rangex][rangey], town t[nb])
{
size_t i,j,k;

fprintf(stderr, "begin_computation_... \n");

for(i=0; i<rangex; i++)
for(j=0; j<rangey; j++) {
pt[i][j].latitude =(xmin+step*i)+180/M_PI;
pt[i][j].longitude =(ymin+step*j)+180/M_PI;
pt[i][j].stock =0.;
for(k=0; k<nb; k++) {
data_t tmp = 6368.* acos(cos(xmin+step*i)+cos( t[k].latitude
+ cos((ymin+step*j)-t[k].longitude)
+ sin(xmin+step*i)*sin(t[k].latitude));
if( tmp < range )
pt[i][j].stock += t[k].stock / (1 + tmp);
}
}
fprintf(stderr, "end_computation_... \n");
}

```

Programming environments for big data processing on modern parallel architectures
HPC Project Ronan Keryell 83 / 101

Hyantes (III)

Example given in par4all.org distribution

Programming environments for big data processing on modern parallel architectures
HPC Project Ronin KerVELL 84 / 101

Hyantes (IV)

OpenMP code:

```
void run(data_t xmin, data_t ymin, data_t xmax, data_t ymax, data_t step, d
{
    size_t i, j, k;
    fprintf(stderr, "begin_computation_...\n");
    #pragma omp parallel for private(k, j)
    for(i = 0; i <= 289; i += 1)
        for(j = 0; j <= 298; j += 1) {
            pt[i][j].latitude = (xmin+step*i)*180/3.14159265358979323846;
            pt[i][j].longitude = (ymin+step*j)*180/3.14159265358979323846;
            pt[i][j].stock = 0;
            for(k = 0; k <= 2877; k += 1) {
                data_t tmp = 6368.*acos(cos(xmin+step*i)*cos(t[k].latitude)*cos
                // (tmp<range)
                pt[i][j].stock += t[k].stock/(1+tmp);
            }
            fprintf(stderr, "end_computation_...\n");
        }
    void display(town pt[290][299])
    {

```

Programming environments for big data processing on modern parallel architectures
HPC Project Ronin KerVELL 85 / 101

Hyantes (V)

```
size_t i, j;
for(i = 0; i <= 289; i += 1) {
    for(j = 0; j <= 298; j += 1)
        printf("%f %f %f\n", pt[i][j].latitude, pt[i][j].longitude, pt
        printf("\n");
}

```

Programming environments for big data processing on modern parallel architectures
HPC Project Ronin KerVELL 86 / 101

Hyantes (VI)

Generated GPU code:

```
void run(data_t xmin, data_t ymin, data_t xmax, data_t ymax, data_t step, d
town pt[290][299], town t[2878])
{
    size_t i, j, k;
    //PIPS generated variable
    town (*P_0)[2878] = (town *) [2878] 0, (*P_1)[290][299] = (town *) [290
    fprintf(stderr, "begin_computation_...\n");
    P4A_accel_malloc(&P_1, sizeof(town[290][299]) - 1 + 1);
    P4A_accel_malloc(&P_0, sizeof(town[2878]) - 1 + 1);
    P4A_copy_to_accel(pt, *P_1, sizeof(town[290][299]) - 1 + 1);
    P4A_copy_to_accel(t, *P_0, sizeof(town[2878]) - 1 + 1);
    p4a_kernel_launcher_0(*P_1, range, step, *P_0, xmin, ymin);
    P4A_copy_from_accel(pt, *P_1, sizeof(town[290][299]) - 1 + 1);
    P4A_accel_free(*P_1);
    P4A_accel_free(*P_0);
    fprintf(stderr, "end_computation_...\n");
}
void p4a_kernel_launcher_0(town pt[290][299], data_t range, data_t step,
data_t xmin, data_t ymin)

```

Programming environments for big data processing on modern parallel architectures
HPC Project Ronin KerVELL 87 / 101

Hyantes (VII)

```

//PIPS generated variable
size_t i, j, k;
P4A_call_accel_kernel_2d(p4a_kernel_wrapper_0, 290, 299, i, j, pt, range,
step, t, xmin, ymin);
}
P4A_accel_kernel_wrapper void p4a_kernel_wrapper_0(size_t i, size_t j, town
data_t range, data_t step, town t[2878], data_t xmin, data_t ymin)
{
    // Index has been replaced by P4A_vp_0:
    i = P4A_vp_0;
    // Index has been replaced by P4A_vp_1:
    j = P4A_vp_1;
    // Loop nest P4A end
    p4a_kernel_0(i, j, &pt[0][0], range, step, &t[0], xmin, ymin);
}
P4A_accel_kernel void p4a_kernel_0(size_t i, size_t j, town *pt, data_t ran
data_t step, town *t, data_t xmin, data_t ymin)
{
    //PIPS generated variable
    size_t k;
    // Loop nest P4A end

```

Programming environments for big data processing on modern parallel architectures
HPC Project Ronin KerVELL 88 / 101

Hyantes (VIII)

```


if (i < 289 && j < 298) {
    pt[289+i][j].latitude = (xmin+step*i)*180/3.14159265358979323846;
    pt[289+i][j].longitude = (ymin+step*j)*180/3.14159265358979323846;
    pt[289+i][j].stock = 0;
    for(k = 0; k <= 2877; k += 1) {
        data_t tmp = 6368.*acos(cos(xmin+step*i)*cos((t+k).latitude)*co
        -(t+k).longitude)*sin(xmin+step*i)*sin((t+k).latitude)
        // (tmp<range)
        pt[289+i][j].stock += t[k].stock/(1+tmp);
    }
}
}
}

```

Programming environments for big data processing on modern parallel architectures
HPC Project Ronin KerVELL 89 / 101

Stars-PM

- Particle-Mesh N-body cosmological simulation
- C code from Observatoire Astronomique de Strasbourg
- Use FFT 3D
- Example given in par4all.org distribution




Programming environments for big data processing on modern parallel architectures
HPC Project Ronan Keryell 90 / 101

Stars-PM time step

```

void iteration(coord pos[NP][NP][NP],
              coord vel[NP][NP][NP],
              float dens[NP][NP][NP],
              int data[NP][NP][NP],
              int histo[NP][NP][NP]) {
    /* Split space into regular 3D grid: */
    discretisation(pos, data);
    /* Compute density on the grid: */
    histogram(data, histo);
    /* Compute attraction potential
    in Fourier's space: */
    potential(histo, dens);
    /* Compute in each dimension the resulting forces and
    integrate the acceleration to update the speeds: */
    forcex(dens, force);
    updatevel(vel, force, data, 0, dt);
    forcey(dens, force);
    updatevel(vel, force, data, 1, dt);
    forcez(dens, force);
    updatevel(vel, force, data, 2, dt);
    /* Move the particles: */
    updatepos(pos, vel);
}
    
```




Programming environments for big data processing on modern parallel architectures
HPC Project Ronan Keryell 91 / 101

Stars-PM & Jacobi results with p4a 1.1.2

- 2 Xeon Nehalem X5670 (12 cores @ 2,93 GHz)
- 1 GPU nVidia Tesla C2050 CUDA 3.2
- Automatic call to CuFFT instead of FFTW (stubs...)
- 150 iterations of Stars-PM

Execution time	p4a	Simulation Cosmo.			Jacobi
		32 ³	64 ³	128 ³	
Sequential	(gcc -O3)	0,68	6,30	98,4	24,5
OpenMP 6 threads	--openmp	0,16	1,28	16,6	13,8
CUDA base	--cuda	0,88	5,21	31,4	67,7
Optim. comm. 1.1	--cuda --com-opt.	0,20	1,17	8,9	6,5
Reduction Optim. 1.1.2	--cuda --com-opt.	0,10	0,32	2,1	3,8
Manual optim.	(gcc -O3)	0,05	0,26	1,8	


p4a 1.1.2 introduce generation of CUDA atomic updates for PIPS detected reductions. Other solution to investigate: CuDPP call generation



Programming environments for big data processing on modern parallel architectures
HPC Project Ronan Keryell 92 / 101

Outline

- 1 Hardware architectures
 - Classical multicores
 - GPU
 - MP-SoC underworlds
- 2 Software environments
 - Programming challenges
 - Multicores
 - GPU
 - Application libraries
- 3 Par4All
 - GPU code generation
 - Code generation for SCMP
- 4 Conclusion




Programming environments for big data processing on modern parallel architectures
HPC Project Ronan Keryell 93 / 101

SCMP computer

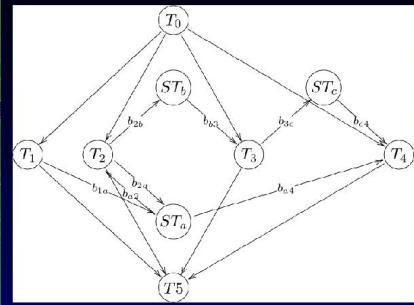
Future revolutions
Software radio, cognitive radio, passive radar, compressed sensing...

- Embedded accelerator developed at French CEA
 - ▶ Task graph oriented parallel multiprocessor
 - ▶ Hardware task graph scheduler
 - ▶ Synchronizations
 - ▶ Communication through memory page sharing
- Generating code from THALES (TCF) GSM sensing application in SCALOPES European project
- Reuse output of PIPS GPU phases + specific phases
 - ▶ SCMP code with tasks
 - ▶ SCMP task descriptor files
- Adapted Par4All Accel run-time




Programming environments for big data processing on modern parallel architectures
HPC Project Ronan Keryell 94 / 101

SCMP tasks



In general case, different tasks can produce data in unpredictable way: use helper data server tasks to deal with coherency when several producers



Programming environments for big data processing on modern parallel architectures
HPC Project Ronan Keryell 95 / 101

SCMP task code (before/after)

```

int main() {
    int i, t, a[20], b[20];
    for (t=0; t < 100; t++)
    {
        kernel_tasks_1:
        for(i=0; i < 10; i++)
            a[i] = i+t;
        kernel_tasks_2:
        for(i=10; i < 20; i++)
            a[i] = 2*i+t;
        kernel_tasks_3:
        for(i=10; i < 20; i++)
            printf("a[%d]=_%d\n",
                i, a[i]);
    }
    return (0);
}

int main() {
    P4A_scmp_reset();
    int i, t, a[20], b[20];
    for(t = 0; t <= 99; t += 1) {
        [...]
        //PIPS generated variable
        int (*P4A_a__1)[10] = (int (*)[10]) 0;
        P4A_scmp_malloc((void **) &P4A_a__1,
            sizeof(int)*10, P4A_a__1_id,
            P4A_a__1_prod_p || P4A_a__1_cons_p, P4A__
            (scmp_task_2_p)
            for (i = 10; i <= 19; i += 1)
                (*P4A_a__1)[i-10] = 2*i+t;
        P4A_copy_from_accel_id(sizeof(int), 20, 10,
            P4A_sam_server_a_p ? ka[0] : NULL, *P4A__
            P4A_a__1_id, P4A_a__1_prod_p || P4A_a__
            P4A_scmp_dealloc(P4A_a__1, P4A_a__1_id,
                P4A_a__1_prod_p || P4A_a__1_cons_p, P4A
        [...]
    }
    return(ev_T004);
}
    
```

Programming environments for big data processing on modern parallel architectures
HPC Project Ronin KERVELL 98 / 101

Performance of GSM sensing on SCMP

- Speed-up on 4 PE SCMP:
 - ×2.35 with manual parallelization by SCMP team
 - ×1.86 with automatic Par4All parallelization
- Still big memory overhead
- To optimize...
 - Use these techniques to generate automate cloud-ification ☺

Programming environments for big data processing on modern parallel architectures
HPC Project Ronin KERVELL 97 / 101

Performance of GSM sensing on SCMP

- Speed-up on 4 PE SCMP:
 - ×2.35 with manual parallelization by SCMP team
 - ×1.86 with automatic Par4All parallelization
- Still big memory overhead
- To optimize...
 - Use these techniques to generate automate cloud-ification ☺

Programming environments for big data processing on modern parallel architectures
HPC Project Ronin KERVELL 97 / 101

Outline

- Hardware architectures
 - Classical multicores
 - GPU
 - MP-SoC underworlds
- Software environments
 - Programming challenges
 - Multicores
 - GPU
 - Application libraries
- Par4All
 - GPU code generation
 - Code generation for SCMP
- Conclusion

Programming environments for big data processing on modern parallel architectures
HPC Project Ronin KERVELL 98 / 101

French advantage: Saint (Holly) Cloud localization

Programming environments for big data processing on modern parallel architectures
HPC Project Ronin KERVELL 99 / 101

Saint Cloud gatekeeper & massive virtual I/O

Programming environments for big data processing on modern parallel architectures
HPC Project Ronin KERVELL 100 / 101

Conclusion

Conclusion

- Quantum change in data amount to process
- Getting out of old tracks for innovative solutions
- Efficiency → future will be more and more heterogeneous ☺
 - Processors everywhere
 - Smart routers with on-the-fly processing
- Hardware is no longer oblivious → Challenge of plain object oriented modeling ☺
- Low latency is difficult (high-frequency trading...)
- Need expertise in hardware, software & applications...
- No unique software environment ☺
- Opportunity to modernize legacy applications
- Good trade-off between efficiency and portability in some areas: OpenCL
- Automatic tools such as Par4All can reduce time-to-market to generate // code to a given target (from low- to high-level)

Programming environments for big data processing on modern parallel architectures
 HPC Project Ronan Keryell 101 / 101

Table of content

- HyperParallel Technologies (1992-1998) 2
- HyperParallel Technologies (1992-1998) 3
- Present motivations: reinterpreting Moore's law 4
- Heterogeneous parallelism 4
- The "Software Crisis" 7
- Time to be back in parallelism! 8
- HPC Project hardware: WildNode from Wild Systems 9
- HPC Project software and services 10
- Efficient big data architectures 11
- Hardware architectures 12
 - Classical multicores 13
 - Intel Nehalem 14
 - AMD Opteron 8180 (2011) 18
 - IBM Power 7 (2010) 19
 - GPU 20
 - Off-the-shelf AMD/ATI Radeon HD 6970 GPU 21
 - Radeon HD 6870 — big picture 22
 - Off-the-shelf NVIDIA Tesla Fermi M2090 & GTX580 23
 - GF100 Stream Multiprocessor 24
 - MP-Soc underworlds 25
 - ARM yourself 26
 - Tilera TilePro64 28
 - MPPA de Kalray (the French touch) 30
 - FPGA 32
 - Convey HPC-10x 34
 - Anton computer from D.E. Shaw 35
 - ProtonShader 36
 - OpenFlow 37
- Software environments 38
 - Programming challenges 39
 - Parallel application owners 40
 - Extracting parallelism in applications... 41
- ... but multidimensional heterogeneity! 42
 - Multicores 43
 - Outline 43
 - OpenMP 44
 - Modèles d'exécution d'OpenMP 45
 - Exemple 47
 - Task on OpenMP 3.0 48
 - Multimedia SIMD extensions 49
 - Message Passing Interface (MPI) 51
 - Thread Building Blocks (TBB) 53
 - GPU 55
 - Basic GPU programming model 56
 - GPU execution model 57
 - From Hardware constraints to programming style 58
 - Programming CUDA 60
 - OpenCL 62
 - CUDA or OpenCL? 66
 - Application libraries 67
 - Outline 68
 - Bibliothèques mathématiques 68
 - Par4All 69
 - Use the Source, Luke... 70
 - Not reinventing the wheel... No NIH syndrome please! 71
 - Current PPS usage 72
 - Par4All usage 73
 - Par4All == PyPS scripting in the backstage 74
 - Coding rules 77
 - Parallelization to OpenMP 78
 - OpenMP output sample 79
 - GPU code generation 80
 - Outline 81
 - Basic GPU execution model 82
 - SciLab & Matlab 83
 - Hyattos 83
 - Stars-PM 91
 - Stars-PM time step 92

Programming environments for big data processing on modern parallel architectures
 HPC Project Ronan Keryell 101 / 101

Table of content

- Stars-PM & Jacobi results with p4s 1.1..2 93
- Performance of GSM sensing on SCMP 99
- Code generation for SCMP 100
- Conclusion 100
- Outline 101
- SCMP computer 95
- French advantage: Saint (Holly) Cloud localization 101
- SCMP tasks 96
- Saint Cloud gatekeeper & massive virtual I/O 102
- SCMP task code (before/after) 97
- Conclusion 103
- Performance of GSM sensing on SCMP 98
- You are here! 105

Programming environments for big data processing on modern parallel architectures
 HPC Project Ronan Keryell 101 / 101

2.6 Denis Caromel (ActiveEon-INRIA)

Solutions ProActive pour Workflows, Map/Reduce, Matlab/Scilab, CPU/GPU

ProActive Parallel Suite (<http://ProActive.inria.fr>), un projet *Open Source* d'OW2, offre une solution flexible pour regrouper des ressources de calcul et offrir aux entreprises un accès simple et unifié à ces ressources par le biais de Portails et d'API. ProActive optimise l'exécution des applications les plus exigeantes, les *workflows* d'entreprises, les simulations numériques et financières, l'analyse des données (avec un *Map/Reduce* qui supporte les APIs Hadoop sans nécessiter un *cluster* dédié). Un mécanisme de sélection de ressources permet de combiner les exécutions sur CPU et GPUs, sur des tâches simples ou au sein même de *workflows* dynamiques. Des analyses de données parallèles sur *Cloud* peuvent être initiées directement sans quitter les environnements Matlab et Scilab. Des *benchmarks* dans les *biotechs* et des démonstrations interactives de *Map/Reduce* seront présentées sur une plate-forme en production.



Solutions ProActive pour Workflows Map/Reduce, Matlab/Scilab, CPU/GPU

Denis Caromel (INRIA & ActiveEon)
Cédric Dalmasso (ActiveEon)

Accelerate and Orchestrate Enterprise Applications

Hybrid Cloud Solutions (CPU+GPU, Private, Public Burst, Multi-Tenants)

Le déluge de données, Ecole Polytechnique, Palaiseau, June 9th 2011



□ Researchers (5):

- D. Caromel (UNSA, Det. INRIA)
- E. Madelaine (INRIA)
- F. Baude (UNSA)
- F. Huet (UNSA)
- L. Henrio (CNRS)

□ PhDs (11):

- Antonio Cansado (INRIA, Conic)
- Brian Amedro (SCS-Agos)
- Cristian Ruz (INRIA, Conicyt)
- Elton Mathias (INRIA-Cordi)
- Imen Filali (SCS-Agos / FP7 SO)
- Marcela Rivera (INRIA, Conicyt)
- Muhammad Khan (STIC-Asia)
- Paul Naoumenko (INRIA/Région)
- Viet Dung Doan (FP6 Bionets)
- Virginie Contes (SOA4ALL)
- Guilherme Pezzi (AGOS, CIPRA)

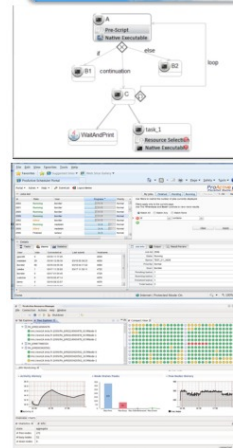
□ + Visitors + Interns



Located in Sophia Antipolis, between
Nice and Cannes,
Visitors Welcome!

ActiveEon Overview

- **ActiveEon**, a software company born of INRIA, founded in 2007HQ in the French scientific park Sophia Antipolis
- **Co developing** with INRIA *ProActive Parallel Suite*®, a Professional Open Source middleware for parallel, distributed, multi-core computing 30 peoples in total
- Core **mission**: Scale Beyond Limits
- Providing a **full range of services** for ProActive Parallel Suite
- **Worldwide** customers and production users:



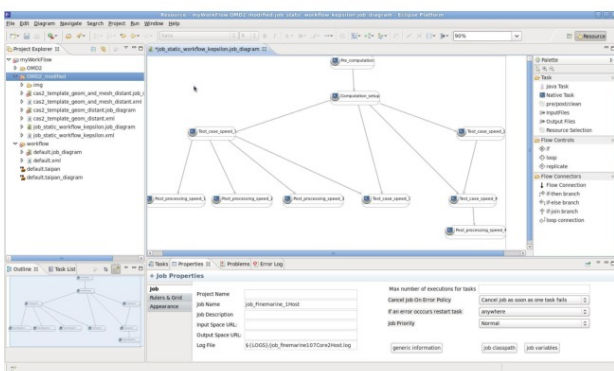
Workflow Execution
Studio Editor and Visualization

Portal, Multi-Application & Multi-Tenant
Enterprise Orchestration

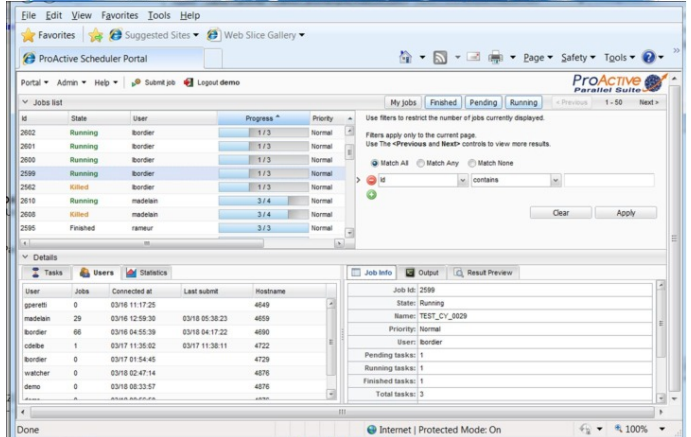
Physical and Virtual Machines Management

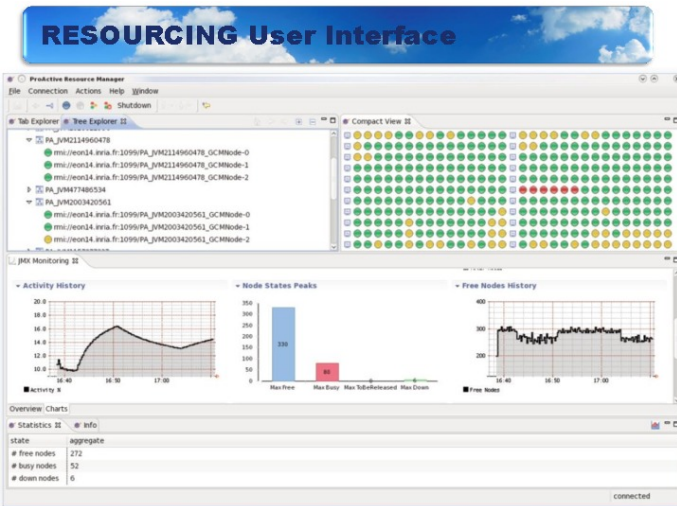


Workflow Studio



ProActive Orchestration Portal





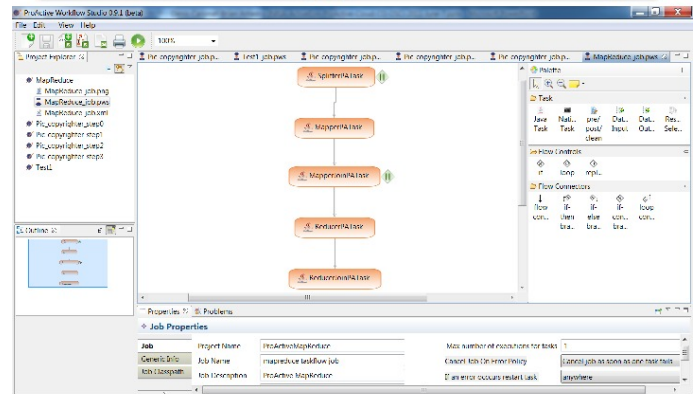
ProActive Parallel Suite Map Reduce

Live Demo

ProActive MapReduce (CO, SP2, Task 2.1)

- ❑ Same APIs as Hadoop (Easy switch from Hadoop to ProActive)
- ❑ Does not requires an HDFS File System
- ❑ Runs on general purpose, Multi-tenant, Multi-Applications Grids and Clouds
- ❑ Available as PaaS in Java

Workflow ProActive MapReduce



ProActive MapReduce vs. Hadoop+HDFS

File Size	Sequential	Hadoop	PA MapReduce	Speedup
0.7 GB	5m 04s	1m 17s	1m 05s	4.6
4.3 GB	25m 31s	2m 30s	2m 20s	10.9
7.3 GB	46m 00s	3m 31s	3m 30s	13.1
20 GB	2h 07m 00s	8m 30s	7m 09s	17.8
50 GB	5h 19m 00s	21m 05s	25m 11s	12.7
100 GB	10h 38m 00s	43m 23s	58m 42s	10.9

- ❑ Data available in a NAS (General purpose storage)
- ❑ Transfer to HDFS for Hadoop
- ❑ Used directly without copy for ProActive
- ❑ Use Case of Map/Reduce on fresh data
- ❑ Different ProActive Map/Reduce configuration for recurrent MR on in place Data (e.g. ProActive HDFS interface)



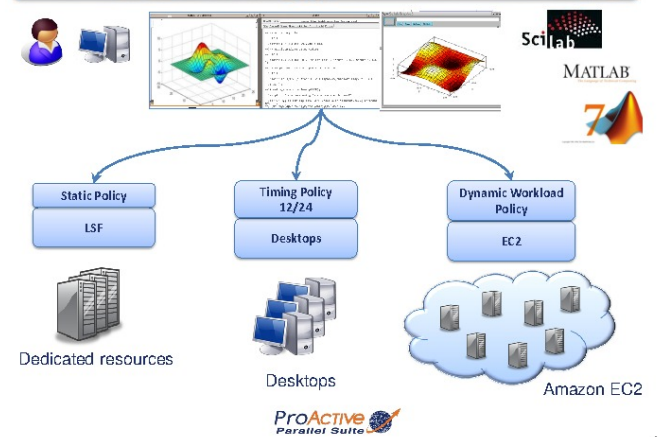


13



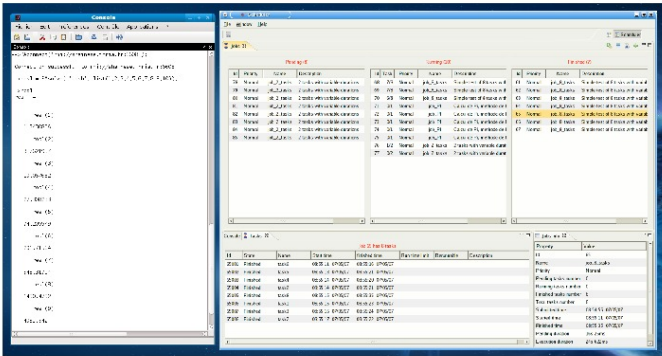
13

Integration with Scilab and Matlab



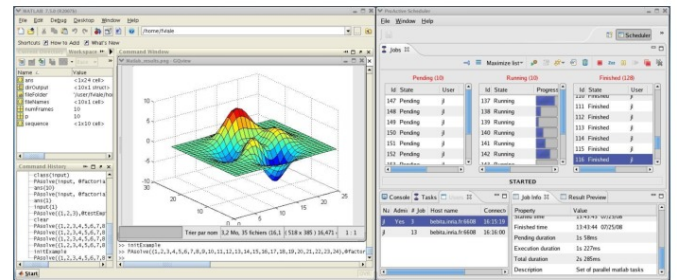
14

Interface ProActive ↔ Scilab



15

Interface ProActive ↔ Matlab



16



Live Demo

17



17

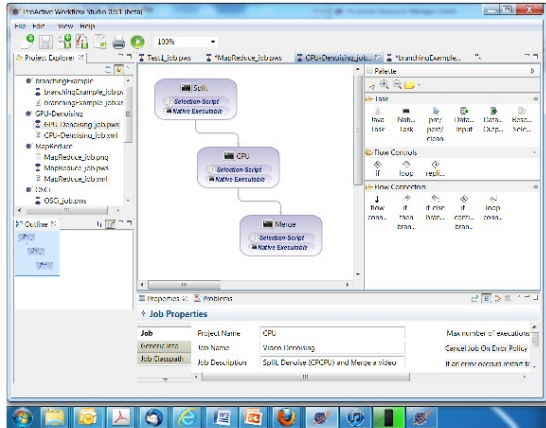
CPU + GPU ProActive Workflows

- Resource selection for each Task of a ProActive Workflow
- Selection of Host with GPU capacity
- Data Transfer to the GPU Host
- Configuration of GPU Capacity at the level of Admin (Number of GPU Nodes, size)
- Freedom to request one or several GPU capacities for one GPU program
- Global Scheduling (Multi-Tenant, Multi-Application) of GPU Tasks

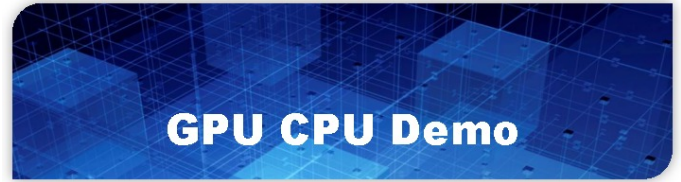


18

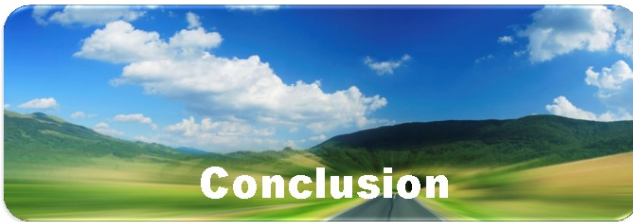
Workflow ProActive for CPU and GPU



19



20



21

Conclusion:



Workflow Studio editor
Workflow Execution (with Visualization)
Map/Reduce, Matlab, Scilab, CPU and GPU

Portal, and APIs (Java, REST, ...)
Multi-Application & Multi-Tenant
Enterprise Orchestration

Physical and Virtual Machines Management
(Hyper-V, VMware, VirtualBox, KVM, Qemu, Xen)
Public Cloud (EC2, Windows Azure), ...



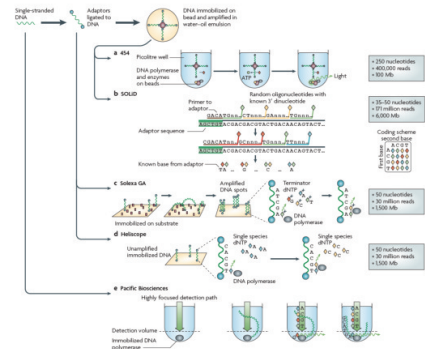
22

2.7 Nicolas Pons (INRA)

La métagénomique, un défi supplémentaire pour la loi de Moore

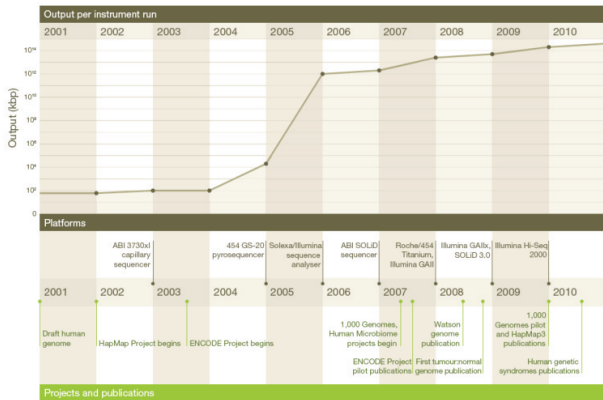
Depuis le séquençage du premier génome en 1995, la production de données de séquençage d'ADN a révolutionné les possibilités de compréhension du vivant par la biologie moléculaire. Avec l'arrivée des technologies de séquençage à très haut-débit, on assiste aujourd'hui à une explosion des volumes de données avec un doublement des bases de données de séquence tous les 6 mois et une augmentation du débit d'acquisition d'un facteur 1000. Ce déluge de données ouvre de nouvelles perspectives scientifiques notamment dans le domaine de la métagénomique qui vise à caractériser l'ensemble des génomes bactériens d'un écosystème complexe : il est désormais possible de quantifier les génomes, gènes et fonctions de ces écosystèmes. Le traitement de ces *big data* constitue un défi majeur tant en matière d'optimisation des calculs qu'en matière de stockage et de leur mise à disposition aux biologistes. Nous illustrerons ces défis à travers l'exemple des projets `MetaHIT` et `MicroObes` qui proposent d'étudier le génome de l'ensemble des bactéries constituant la flore intestinale humaine afin de caractériser ses fonctions et ses implications sur la santé.

Next Generation Sequencing



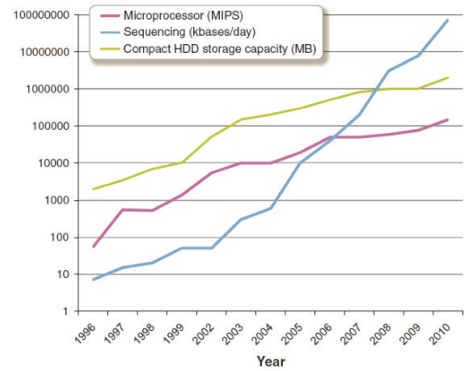
Different sequencing instruments

Instrument	Run time ^a	Millions of reads/run	Bases/read ^b	Yield Mb/run	Reagent cost/run ^c	Reagent cost/Mb	Minimum unit cost (% run) ^d
3730xl (capillary)	2 h	0.000096	650	0.06	\$96	\$1500	\$6 (1%)
Ion Torrent - '314' chip	2 h	0.10	100	>10	\$500	<\$50	~\$750 (100%)
454 GS Jr. Titanium	10 h	0.10	400	50	\$1100	\$22	\$1500 (100%)
Starlight [®]	†	>0.01	>1000	†	†	†	†
PaSeq RS	0.5-2 h	0.01	860-1100	5-10	\$110-900	\$11-180	†
454 FLX Titanium	10 h	1	400	500	\$6200	\$12.4	\$2000 (10%)
454 FLX+	18-20 h	1	700	900	\$6200	\$7	\$2000 (10%)
Ion Torrent - '316' chip ^a	2 h	1	>100	>100	\$750	<\$7.5	~\$1000 (100%)
Helicos ¹	N/A	800	35	28 000	N/A	N/A	\$1100 (2%)
Ion Torrent - '318' chip ^a	2 h	4-8	>100	>1000	~\$925	~\$0.93	~\$1200 (100%)
Illumina MiSeq ¹	26 h	3.4	150 + 150	1020	\$750	\$0.74	~\$1000 (100%)
Illumina ScanSQ	8 days	250	100 + 100	50 000	\$10 220	\$0.20	\$3000 (14%)
Illumina GAIIx	14 days	320	150 + 150	96 000	\$11 524	\$0.12	\$3200 (14%)
SOLID - 4	12 days	>840 ^b	50 + 35	71 400	\$8128	<\$0.11	\$2500 (12%)
Illumina HiSeq 1000	8 days	500	100 + 100	100 000	\$10 220	\$0.10	\$3000 (12%)
Illumina HiSeq 2000	8 days	1000	100 + 100	200 000	\$20 120 ^b	\$0.10	\$3000 (6%)
SOLID - 5500 (PT) ^a	8 days	>700 ^b	75 + 35	77 000	\$6101	<\$0.08	\$2000 (12%)
SOLID - 5500xl (4hq) ^a	8 days	>1410 ^b	75 + 35	155 100	\$10 503 ^b	<\$0.07	\$2000 (12%)
Illumina HiSeq 2000 - v3 ^a	10 days	≤3000	100 + 100	≤600 000	\$23 470 ^b	≥\$0.04	~\$3500 (6%)



Sequencing Progress vs Compute and Storage

Moore's and Kryder's Laws fall far behind



From kilobytes to terabytes !!!




Year	Sequenced object	Sequencing unit
1990	Genes & operons	Kbase
1995	Bacterial genome	Mbase
2001	Human genome	Gbase
2010	Human metagenome	Tbase

Amount of sequence we generate has increased 10^9 times in 20 years, greatly exceeding the Moore's law







New bioinformatics challenges

Data management



Data analysis



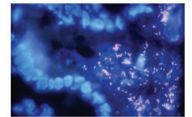
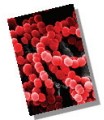
Data browsing

METEOR around iMOMi framework



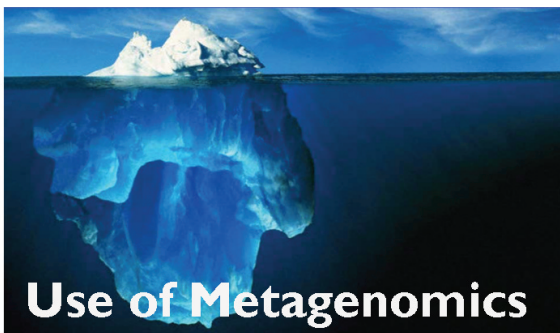
The human intestinal microbiota is a forgotten organ...

- ✓ 100 trillion microorganisms ; 10-fold more cells than the human body; 2 kg of mass!
- ✓ Interface between food and epithelium
- ✓ In contact with the 1st pool of immune cells and the 2nd pool of neural cells of the body



...with a major role in health & disease !

Most of microorganisms are unknown and uncultivable...



Use of Metagenomics

What is metagenomics ?

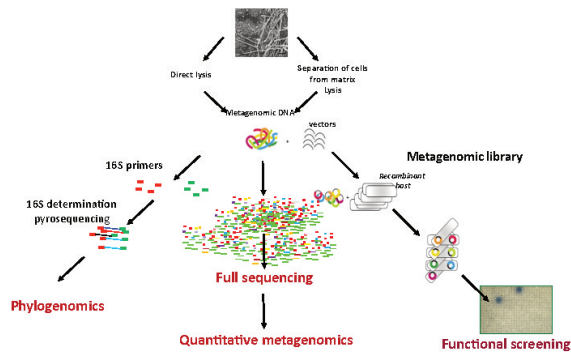
Metagenome

can be defined as the ensemble of genes of the microbes from a given ecological niche.

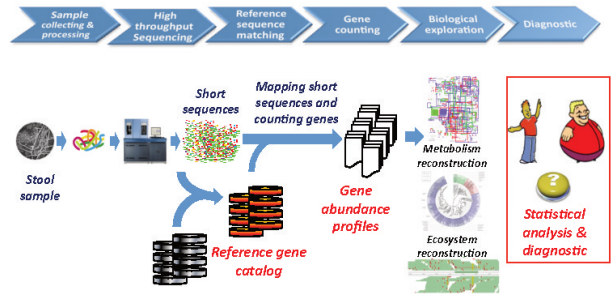
Metagenomics

allows to characterize composition, properties and dynamics of a microbiome by studying the metagenome.

Different types of metagenomics



Quantitative metagenomics pipeline



A powerful microscope!

The MetaHIT project (Dr. S.D. Ehrlich)



- ✓ Establishing a **reference gene catalog** by metagenomic & genomic sequencing of the Human GI tract microbes
- ✓ Developing generic **tools for profiling** the GI tract microbiota genes : arrays and high throughput DNA sequencing

Using Quantitative Metagenomics to **search associations of microbial genes and chronic disease** in Obesity and Inflammatory Bowel Diseases

Illumina sequencing

Samples	124 individuals (85 Danes, 39 Spaniards)	
Library type	15 samples	200bp
	109 samples	140bp 350bp
Sequencing type	Paired-end (PE) sequencing	
Read length (bp)	45 b (15 samples) 75 b (109 samples)	
Tags per sample	31million ±0.5 million	

In total, ~0.58 Terabase sequence

Wang Jun et al.

Our other genome : the human intestinal metagenome



March 2010

3.3 million bacterial gene catalog : 150-fold human genome
 Each individual has ~ 540000 of the 3.3 million genes
 85% of abundant gut genes from a cohort of 124 individuals
 70-86% of genes from the US & Japanese studies

A complementary study : Micro-Obes (Dr. J. Doré)

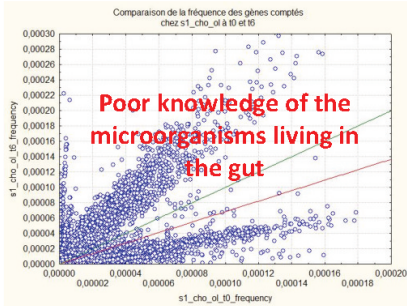


Follow the microbiota dynamic and identify signatures of obesity and nutritional transition in the intestinal microbiote

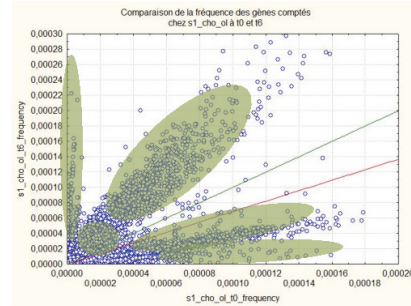
- DNA sequencing (SOLID) of 215 faecal samples of 49 obese individuals sequenced at 3 differents time-points (t_0 , w_6 , w_{12} and extra-time)
 → 300 Gbases sequenced
- Read projection against the MetaHIT gene catalog (Qin et al., 2010)
 (3.3 millions genes)
 → 150 Gbases projected



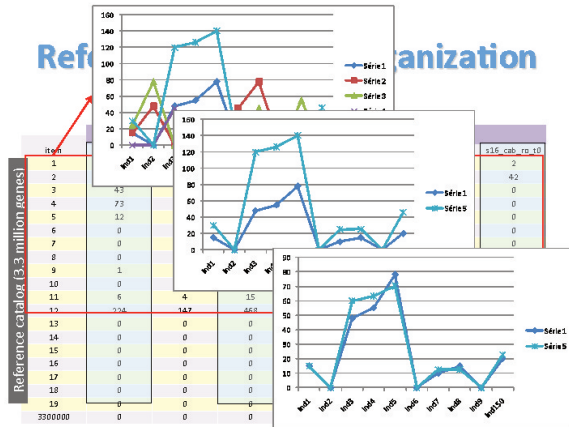
Dynamic of the intestinal microbiota



Dynamic of the intestinal microbiota



Reference catalog reorganization



Reference catalog reorganization

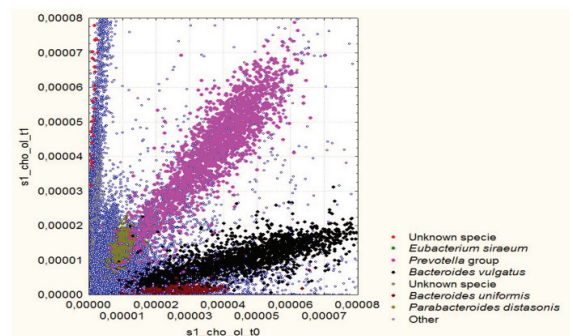
- Hierarchical descendant graph & DAPC clustering (Almeida et al., 2011)
 - By computation of spearman correlation for each couple of gene profile
 - $3.3^6 \times 363 \rightarrow 5^{12}$ correlations to calculate
 - With one CPU : **more than a year to do it...**
- MetaProf (Boumezbeur, Arslan et al., 2011)
 - CUDA programming
 - 1H10 in 40 GPU



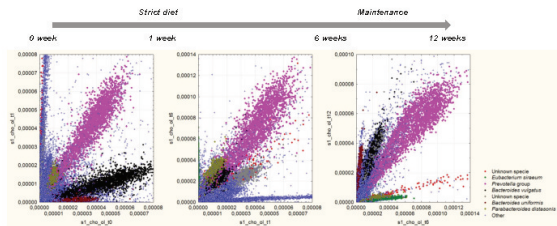
Reference catalog reorganization

- 17317 clusters from 10 to 59353 genes (~2.34M genes, ~71%)
- 628 meta-species with more than 1000 genes

Dynamic of the intestinal microbiota

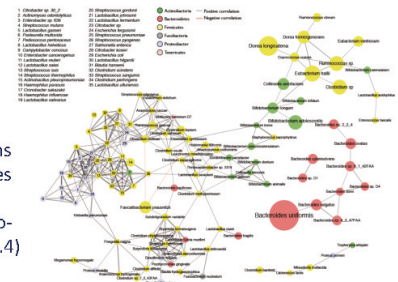


Dynamic of the intestinal microbiota



Important variation of gene composition during the first part of the diet
More stable composition in the second part (stabilization)

Co-variation of bacterial species



Network of relations between the species (circles represent species, lines the co-variation, with $R \geq 0.4$)

Particular constellations of bacterial species in individuals

DTU, Brunak et al.

ARTICLE

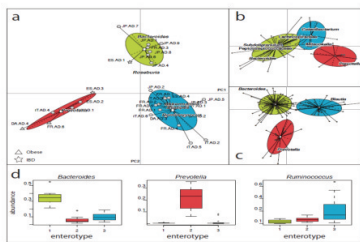
doi:10.1038/nature09944

Enterotypes of the human gut microbiome

Mamunohoyan Arunugam¹, Jeroen Raaijmakers^{1,2}, Eric Pelletier^{3,4,5}, Denis Le Paslier^{6,7,8}, Takaji Yamada⁹, Daniel R. Mende¹, Gabriel R. Ferraz¹⁰, Julien Tarr¹¹, Thomas Brink^{12,13}, Jean-Michel Hainey¹⁴, Marco Bertalan¹⁵, Natalia Borczyk¹⁶, Francisca Casellas¹⁷, Leyden Fernandez¹⁸, Laurent Gaillier¹⁹, Torben Hansen^{20,21}, Masahito Hattori²², Tetuya Hayashi²³, Michiel Kleverbergen²⁴, Kenji Kurakawa²⁵, Martin Leclercq²⁶, Pierrick Levenez²⁷, Chayaporn Manichanh²⁸, H. Bjørn Nielsen²⁹, Trine Nielsen³⁰, Nicolas Pons³¹, Julie Prezelain³², Junjie Qian³³, Thomas Sieberitz³⁴, Sebastian Timm³⁵, David Torrent^{36,37}, Edgardo Ugarte³⁸, Erwin C. Zosenda³⁹, Jun Wang⁴⁰, Francisco Guarner^{41,42,43}, Willem M. de Vos⁴⁴, Søren Brunak⁴⁵, Joel Dor⁴⁶, MetaHT Consortium^{1,47}, S. Dusko Ehrlich⁴⁸ & Peer Bork⁴⁹

May 2011

Europeans, Americans, Asians. n=33; Sanger



Danes n=85; Illumina
US n=154; 454

Enterotypes can be likened to blood groups but the reasons for their existence remains to be elucidated

They should allow patient stratification & aid to develop personalized medicine and nutrition

Where do these studies lead to and when?

- **Diagnostic & prognostic tests – soon**
 - arrays, sequencing, Q-PCR; immunomarkers (?)
- **Better treatments – next**
 - personalized medicine
- **Novel treatments – last**
 - modulation of microbiota
 - Promoters
 - Inhibitors
 - transplantation of microbiota

Dr. S.D. Ehrlich :

“Take-home message:

- **Our other genome** has much more variability than the first one
- **Personalized medicine** should target it”

Conclusion

- **Data deluge is not finished**
 - 3rd generation sequencing
 - In our lab: 2 SOLiD 5500xl sequencers the next week
- **Data managing**
 - dCache / HDF5
 - Cloud
 - International repository
- **Data analysis**
 - GPU oriented programming (OpenGPU project)
- **Data browsing**

BioRad

MetaQuant
 Dusko Ehrlich
 Sean Kennedy
 Nicolas Pons
 Nathalie Galleron
 Benoit Quinquis

Plateforme
 Meta Quant
 HighThroughput Quantitative

Meta HIT

ANR
 Micro-Obes
 Food-Microbiomes

Team « Informatique »
 Jean-Michel Batto
 Pierre Léonard
 Bouziane Moumen

Team « Bactéries Alimentaires et Commensales »

Bioinformatic
 Emmanuelle Le Chatelier
 Mathieu Almeida
 Fouad Boumebeur

Biology
 Christine Delorme
 Eric Guédon
 Séverine Layec
 Ghalia Kaci
 Céline Gautier
 Nicolas Sanchez

Pierre Renault

OpenGPU

http://www.netvibes.com/metahit/Live_News
<http://t.me/metagenomics>
<http://paperkit.net/metahit/microbiomics>

paper 31

2.8 Patrick Fuhmann (DESY-Hamburg)

dCache : scaling out affordable storage.

The presentation will briefly walk through the various facets of the dCache storage technology and its supporting collaborations. Functional objectives will be discussed, as well as some bits and pieces of the technical implementation. The presenter will touch upon the results of an ongoing detailed evaluation of supported file access protocols at the DESY Grid-Lab facility including a discussing on application level behaviour and pitfalls. Finally the most prominent dCache deployments will be presented and the involvement of dCache in other projects.



DCACHE.ORG
DCACHE.ORG
DCACHE.ORG

DCACHE, LARGE SCALE-OUT AFFORDABLE STORAGE

PATRICK FUHRMANN
ON BEHAVE OF THE TEAM



9/5/11

SÉMINAIRE ARISTOTE : BIG DATA, ECOLE POLYTECHNIQUE,
PATRICK FUHRMANN

1



DCACHE.ORG
DCACHE.ORG
DCACHE.ORG



CONTENT

- PEOPLE AND FUNDING
- DEPLOYMENT
- SUPPORTED PROTOCOLS
 - DATA ACCESS
 - STORAGE CONTROL
- MANAGED STORAGE
- DESIGN FACTS

9/5/11

SÉMINAIRE ARISTOTE : BIG DATA, ECOLE POLYTECHNIQUE,
PATRICK FUHRMANN

2



DCACHE.ORG
DCACHE.ORG
DCACHE.ORG

PEOPLE FUNDING

9/5/11

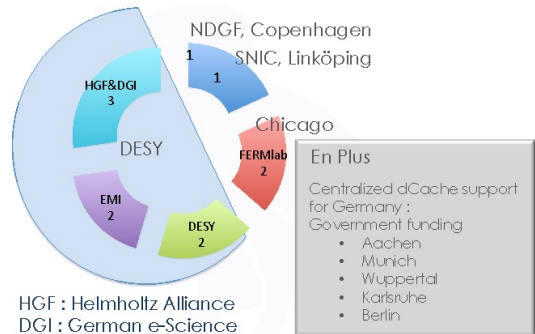
SÉMINAIRE ARISTOTE : BIG DATA, ECOLE POLYTECHNIQUE,
PATRICK FUHRMANN

3



DCACHE.ORG
DCACHE.ORG
DCACHE.ORG

PEOPLE AND FUNDING



HGF : Helmholtz Alliance
DGI : German e-Science
EMI : European Middleware Initiative

9/5/11

SÉMINAIRE ARISTOTE : BIG DATA, ECOLE POLYTECHNIQUE,
PATRICK FUHRMANN

4



DCACHE.ORG
DCACHE.ORG
DCACHE.ORG

TWO WORDS ON EMI

EUROPEAN MIDDLEWARE INITIATIVE

9/5/11

SÉMINAIRE ARISTOTE : BIG DATA, ECOLE POLYTECHNIQUE,
PATRICK FUHRMANN

5



DCACHE.ORG
DCACHE.ORG
DCACHE.ORG

EMI FACTSHEET

EMI Factsheet

- Budget : about 24 Million Euros
- Funding : about 50% by EU-FP7, rest by partners
- Covers : JRA, SA and NA
- Partners : 22
- Middlewares : Arc, gLite, UNICORE and dCache

Logos: Open Science Grid, Google, University of Oslo, Science & Technology, INFN, CINECA, gmet, SWITCH, INET, CERN, KIT, KISTI, etc.

16/09/2010 EMI Overview - EGI TF, Amsterdam
9/5/11 May 25, 2011 SÉMINAIRE ARISTOTE : BIG DATA, ECOLE POLYTECHNIQUE,
PATRICK FUHRMANN

9/5/11

SÉMINAIRE ARISTOTE : BIG DATA, ECOLE POLYTECHNIQUE,
PATRICK FUHRMANN

6



DCACHE.ORG

DEPLOYMENT

- WLCG
- OTHERS

9/5/11

SÉMINAIRE ARISTOTE : BIG DATA, ECOLE POLYTECHNIQUE, PATRICK FUHRMANN

7



WLCG

WLCG : WORLD WIDE LHC COMPUTING GRID
LHC : LARGE HADRON COLLIDER



9/5/11

SÉMINAIRE ARISTOTE : BIG DATA, ECOLE POLYTECHNIQUE, PATRICK FUHRMANN

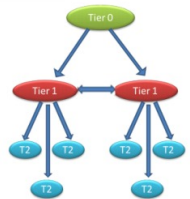
8



WLCG : SPECIFICATION

- Trying to find out why we are heavy (Higgs)
- Maybe more about dark matter and dark energy
- Producing about 15 PBytes per year
- "Raw data" is stored at CERN and at least at another Site (Tier I)

Legend
Tier 0 : CERN
Tier 1 : Counties (11)
Tier 2 : about 200



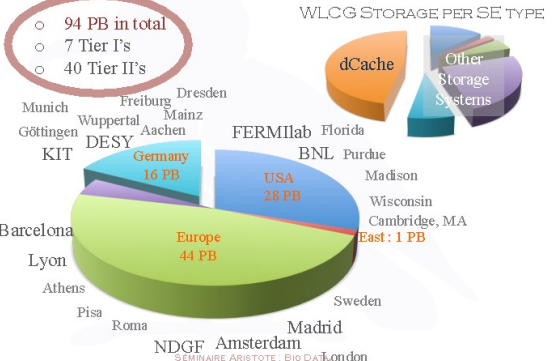
9/5/11

SÉMINAIRE ARISTOTE : BIG DATA, ECOLE POLYTECHNIQUE, PATRICK FUHRMANN

9



STATUS : DCACHE DEPLOYMENT



9/5/11

SÉMINAIRE ARISTOTE : BIG DATA, ECOLE POLYTECHNIQUE, PATRICK FUHRMANN

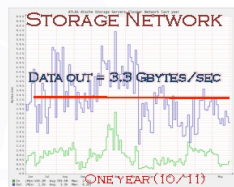
10



LARGEST DCACHE INSTALLATION (BNL) BROOKHAVEN NATIONAL LAB (NEW YORK)

INFORMATION PROVIDED BY HIRONORI ITO (BNL)

- 80 Million files in total
- 85 storage hosts with about 600 pools. (dCache storage unit)
- Total space on disk : 8.8 PBytes (Used are 7 PBytes)
- Total space on tape (HPSS, IBM) : 2.5 PBytes



9/5/11

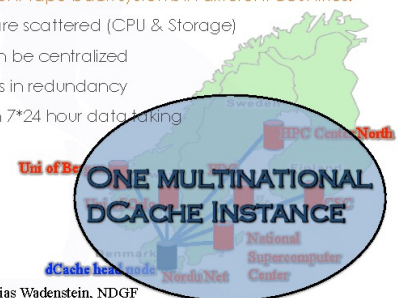
SÉMINAIRE ARISTOTE : BIG DATA, ECOLE POLYTECHNIQUE, PATRICK FUHRMANN

11



MOST INTERESTING DEPLOYMENT

- The 7 biggest Nordic Computer centers form a single Tier I
- Many different tape back systems in different countries.
- Resources are scattered (CPU & Storage)
- Services can be centralized
- Advantages in redundancy
- Especially in 7*24 hour data taking



Slide stolen from Mattias Wadenstein, NDGF

9/5/11

SÉMINAIRE ARISTOTE : BIG DATA, ECOLE POLYTECHNIQUE, PATRICK FUHRMANN

12

OTHER DEPLOYMENTS

DCACHE.ORG
DCACHE.ORG
DCACHE.ORG

Historically

- DESY :
HERA (Zeus, H1)
- FERMILAB :
Tevatron (CDF)
Sloan Digital Sky Survey

NEW DATA INTENSIVE COMMUNITIES

DCACHE.ORG
DCACHE.ORG
DCACHE.ORG



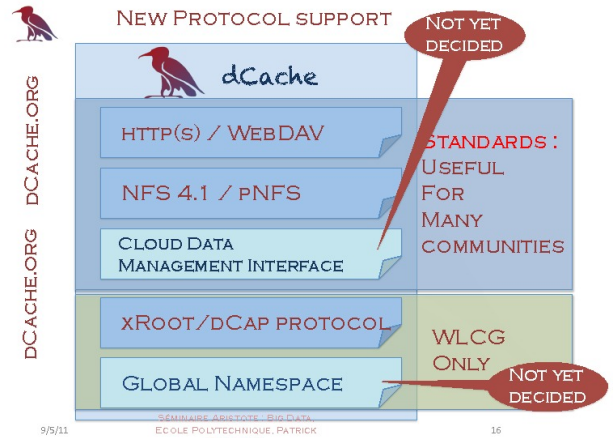
LOFAR
Low Frequency Array
Is using dCache at SARA (Amsterdam) and Jülich (Germany)

CFEL SCIENCE
Center for Free Electron Laser Science
Would like to use dCache at the DESY storage.

SNIC
Swedish National Infrastructure for Computing
Using dCache for Swedish academic purposes.

SUPPORTED PROTOCOLS

DCACHE.ORG
DCACHE.ORG
DCACHE.ORG



PROTOCOL SUPPORT : WEBDAV

DCACHE.ORG
DCACHE.ORG
DCACHE.ORG

- Very useful for new (non-LHC) communities.
- IETF Standard
- Allows "File system like" access with
 - Mac OS
 - Linux
 - Windows



NFS v 4.1 / PNFS

MY FAVORED TOPIC

DCACHE.ORG
DCACHE.ORG
DCACHE.ORG

PROTOCOL SUPPORT : NFSV4.1 /PNFS

center for information technology integration
CITI, at the University of Michigan, is funded by major storage providers to coordinate the pNFS effort and provide reference implementations.

Industry Support - Implementations

- Clients**
 - Linux
 - Sun (Solaris)
- Servers**
 - Desy
 - EMC
 - IBM
 - Linux
 - NetApp
 - Panasas
 - Sun (Solaris)

Group meets three times a year to check interoperability.

Several other implementations have been tested at Bake-a-thons and Connectathons

9/5/11 SEMINAIRE ARISTOTE : BIG DATA, ECOLE POLYTECHNIQUE, PATRICK FUHRMANN 19

PROTOCOL SUPPORT : NFSV4.1 /PNFS

Stolen from : <http://www.pnfs.com/>

9/5/11 SEMINAIRE ARISTOTE : BIG DATA, ECOLE POLYTECHNIQUE, PATRICK FUHRMANN 20

PROTOCOL SUPPORT : NFSV4.1 /PNFS

Stolen from : <http://www.pnfs.com/>

Benefits of Parallel I/O

- ✓ Delivers Very High Application Performance
- ✓ Allows for Massive Scalability without diminished performance

Benefits of NFS (or most any standard)

- Ensures Interoperability among vendor solutions
- Allows Choice of best-of-breed products
- Eliminates Risks of deploying proprietary technology

9/5/11 SEMINAIRE ARISTOTE : BIG DATA, ECOLE POLYTECHNIQUE, PATRICK FUHRMANN 21

PROTOCOL SUPPORT : NFSV4.1 /PNFS

Simplicity

- ✓ Regular mount-point and real POSIX I/O
- ✓ Can be used by unmodified applications (e.g. Mathematica..)
- ✓ Data client provided by the OS vendor
- ✓ Smart caching (block caching) development done by OS vendors

Performance

- ✓ pNFS : parallel NFS (first version of NFS which support multiple data servers)
- ✓ Clever protocols , e.g. Compound Requests

9/5/11 SEMINAIRE ARISTOTE : BIG DATA, ECOLE POLYTECHNIQUE, PATRICK FUHRMANN 22

DCACHE IS MANAGED STORAGE

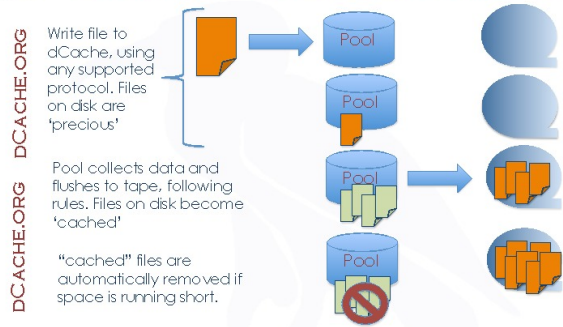
- **MANUAL STORAGE MANAGEMENT**
 - THE STORAGE RESOURCE MANAGER (SRM)
 - STORAGE MIGRATION MODULE
- **AUTOMATIC STORAGE MANAGEMENT**
 - STORAGE ATTRIBUTE BY DIRECTORY
 - HOT SPOT DETECTION
 - RESILIENT MANAGER

9/5/11 SEMINAIRE ARISTOTE : BIG DATA, ECOLE POLYTECHNIQUE, PATRICK FUHRMANN 23

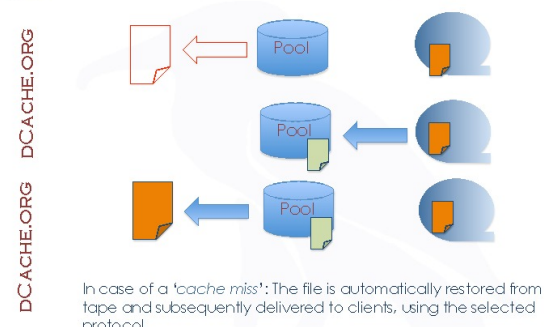
DCACHE IDEA

9/5/11 SEMINAIRE ARISTOTE : BIG DATA, ECOLE POLYTECHNIQUE, PATRICK FUHRMANN 24

NORMAL DCache FILE CYCLE (STORING DATA)

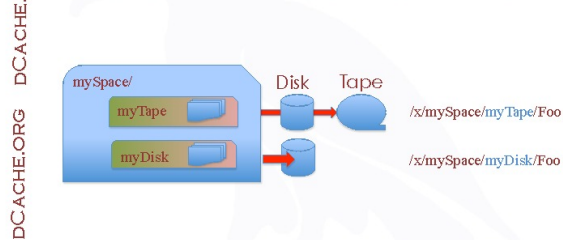


NORMAL DCache FILE CYCLE (RETRIEVING DATA)



MANAGED STORAGE

FINAL DESTINATION DETERMINED BY DIRECTORY



MANAGED STORAGE : SRM

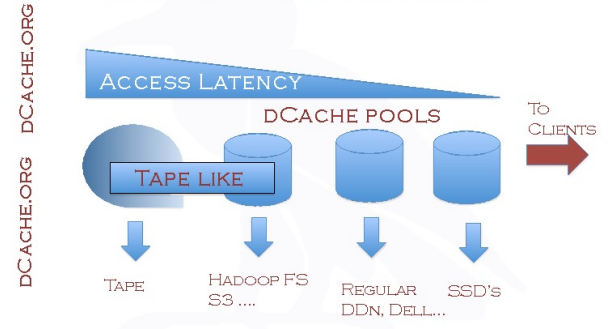
THE STORAGE RESOURCE MANAGER PROTOCOL

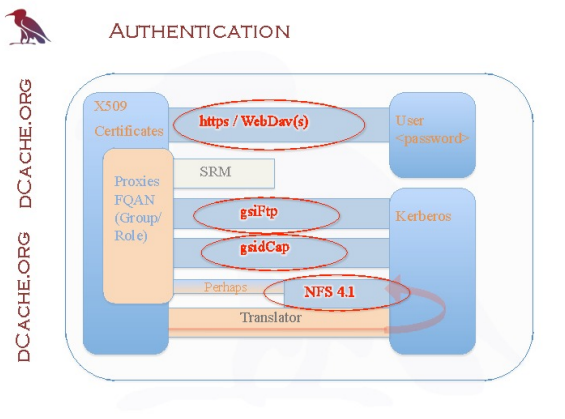
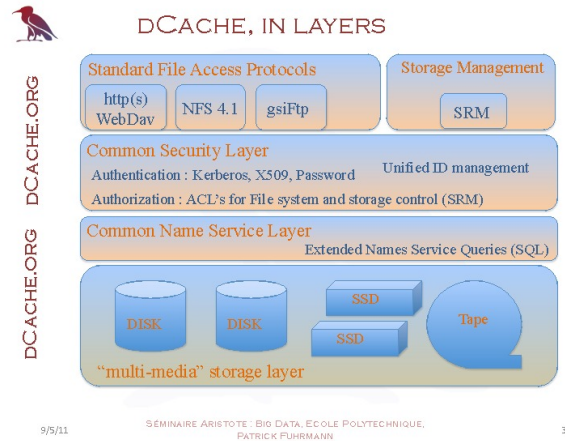
- SRM 2.2 : defined by the Open Grid Forum (OGF)
 - ✓ Defines storage media (Disk/Tape)
 - ✓ Can use "Spaces" (similar to AMAZON buckets) with attributes (disk, tape, size)
 - ✓ Pin / Unpin files
 - ✓ Bring Online file(s), in preparation for a read.
 - ✓ Remote secure protocol with many implementations

MORE MANAGED STORAGE

- Automatic file replication on 'hot spot' detection
 - ✓ If a pool is used heavily, dCache starts to spread files from this pool to other (lazy) pools.
- Resilient manager
 - ✓ On basis of a pool set, a minimum and maximum number of replicas for all files can be defined.
 - ✓ dCache automatically adjusts the replicas if pools go down or are scheduled for maintenance.
- Migration Module
 - ✓ Files can be shuffled around between pools (by rules) to allow to spread load or decommission pools.

DCACHE ALLOWS TO MAKE USE OF DIFFERENT STORAGE CAPABILITIES





- ### NO SECRETS ANYMORE
- > All Java
 - > Name space abstraction
 - > Legacy implementation (PNFS) or
 - > New Implementation (any JDBC DB, def. postgres)
 - > Component communication via message passing:
 - > Private Protocol (Cells) or
 - > Java Messaging Service (JMS)
 - > Scalable components : Protocol Endpoints and data pools
 - > Single Point of failures : namespace and pool/space manager
 - > With 1.9.1.2 (2nd Golden Release) : Very nice configuration system

- ### CONCLUSION
- > dCache is about storing, accessing and managing huge amounts of data.
 - > Depending on the configuration (resilient manager) you may use cheap hardware.
 - > Historically tuned for HEP and WLCG
 - > For about 2 year focusing on more communities, which have been committing themselves to standards (web 2.0)
 - > dCache collaboration nicely distributed amongst Europe and the US.
 - > Funding spread amongst different bodies. (e.g EMI)
 - > More contributions/contributors welcome.


FURTHER READING

WWW.DCACHE.ORG

2.9 Marie-Luce Picard (EDF R&D et ENST-Bilab)

Données massives pour les *smart-grids*


De nombreux projets *smart-grids* voient le jour à travers le monde, motivés par des aspects réglementaires, des contraintes économiques ou la prise en compte de besoins environnementaux ou sociaux. Ces projets reposent sur le déploiement de compteurs communicants et la mise en place d'une infrastructure de communication adéquate. Mais il ne s'agit là que de la première étape de la mutation technique et économique du secteur énergétique. Cette vision long terme de la problématique des réseaux intelligents sous-tend une capacité à gérer et traiter de larges volumes de données, provenant en particulier des compteurs intelligents ou encore de différents capteurs sur le réseau. Dans cette perspective, un certain nombre de travaux ont été menés à EDF R&D et seront présentés dans cet exposé, en particulier : le stockage de grandes quantités de séries temporelles, le traitement temps-réel de courbes de charge imparfaites, les perspectives d'évolution des approches de prévision de consommation en présence de données individuelles massives.



Données massives pour les smart-grids

Marie-Luce PICARD
EDF R&D
marie-luce.picard@edf.fr

9 Juin 2011



CHANGER L'ÉNERGIE ENSEMBLE



Sommaire

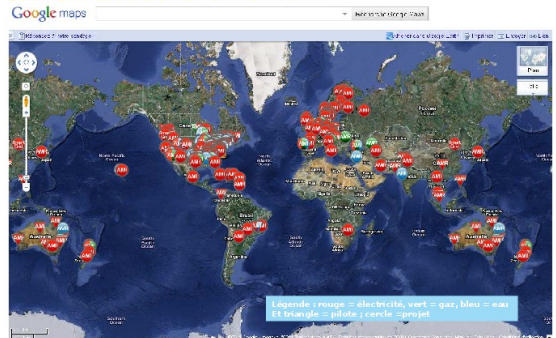
1. Smart-Grids : systèmes électriques intelligents
2. Données et Smart-Grids
3. Travaux menés à EDF R&D
4. Conclusion



Smart Grids : systèmes électriques intelligents



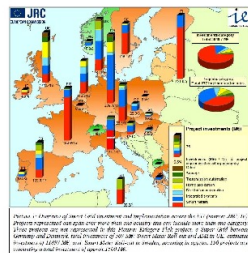
Les réflexions autour des Smart-Grids voient le jour partout dans le monde ...



Un enjeu majeur ... De nombreux investissements



- Les incitations des pouvoirs publics encouragent l'innovation et les expérimentations en la matière (RDG, Ademe, etc.)
- Des investissements importants ont été réalisés par des acteurs TIC ou industriels.
- Des initiatives de mise en service ont été effectuées par de nombreuses Utilités.



Pour plus d'informations sur les investissements dans les Smart Grids, consultez le rapport de l'Agence de l'Énergie (ADEME) intitulé 'Smart Grids : les investissements dans les réseaux électriques intelligents' (juin 2010) disponible sur www.ademe.fr.

Et en particulier en France

Commission de Régulation de l'Énergie (CRE)
EDF : www.amandis.com.fr
linky.edf.fr/distribution.fr

ERDF

Découvrez Linky et 30 services

Bienvenue sur le site ERDF dédié aux compteurs Linky

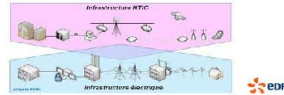
- Décret d'Août 2010, qui indique que :
- À partir du 01/01/2012 tous les nouveaux compteurs doivent être communicants
 - Après le 31/12/2014, 50% des compteurs installés doivent être communicants
 - Après le 31/12/2016, 95% des compteurs installés doivent être communicants.



Smart-Grids : systèmes électriques intelligents

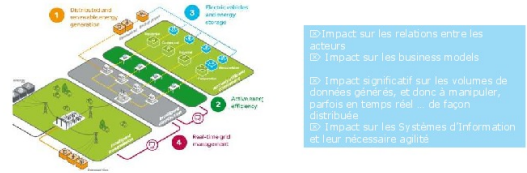
Source – Wikipedia
 A smart grid delivers electricity from suppliers to consumers using digital technology with two-way communications to control appliances at consumers' homes to save energy, reduce cost and increase reliability and transparency. It overlays the electrical grid with an information and networking system, and includes smart meters. Such a modernized electricity network is being promoted by many governments as a way of addressing energy independence, global warming and emergency resilience issues.

- «Un Système électrique intelligent est un système électrique capable d'intégrer de manière intelligente les actions des différents utilisateurs, consommateurs et/ou producteurs afin de maintenir une fourniture d'électricité efficace, durable, économique et sécurisée».
- Passage d'un modèle traditionnel centralisé à un modèle distribué, interconnecté, avec de fortes interactions (avec le client)
- Le réseau se dote de nombreux capteurs et dispositifs de communication bi-directionnel
- Réseau électrique + réseau informatique



Smart-Grids : systèmes électriques intelligents

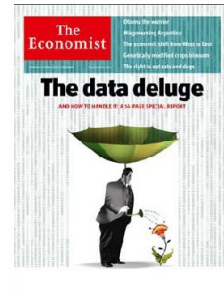
- Les enjeux des Smarts-Grids :
 - Permettre notamment une gestion plus interactive de la demande et de la production décentralisée
 - Favoriser les comportements de MDE (Maîtrise de l'Énergie) : information du client, effacements, demand response ...
 - Gérer les énergies intermittentes et la qualité de fourniture
 - Gérer le développement de nouveaux usages tels que le véhicule électrique



Données et Smart-Grids

EDF R&D : Créer de la valeur et préparer l'avenir

Données massives ?



EDF R&D : Créer de la valeur et préparer l'avenir

Traitement de données massives : quelles technologies ?

Technologies	Exemples
Très grands entrepôts de données	Wal-mart Entrepôts à large échelle : Teradata, Exadata, InfoSphere
Traitement de données à la volée	CME Group Complex Event Processing : StreamBase, IBM InfoSphere Streams
Approches distribuées pour le stockage	Google, Facebook Cloud computing (MapReduce) Hadoop
Data mining à très grande échelle	Amazon Systèmes de recommandation Data mining distribué

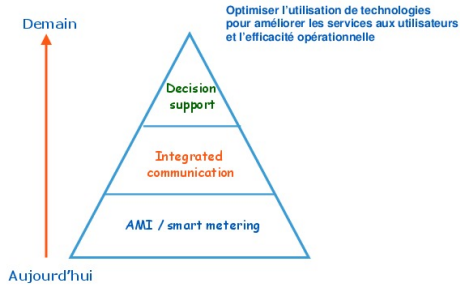
EDF R&D : Créer de la valeur et préparer l'avenir

Qu'en est-il pour les Smart-Grids ?

- La mise en place de fonctionnalités de type Smart-Grids va entraîner une très forte augmentation du volume des données à traiter :
 - Installation des compteurs communicants (35 millions en France)
 - Installation de capteurs divers sur le réseau
 - Évolution des installations chez les clients (objets communicants ...)
 - Peta-octets (dici 5 / 10 ans)
- Ces volumes sont toutefois en-deçà de ceux rencontrés par d'autres secteurs d'activité qui ont d'ores et déjà à leur disposition différents outils technologiques
- Néanmoins on entrevoit la mise en place de traitements complexes :
 - De par la nature des données : des séries temporelles
 - De par leur caractère distribué et de la nécessité de traitements à différentes échelles
 - De par les contraintes temps réel pour certains besoins

A Grid is a Grid: What Does Scale Have to Do with It?(SDG&E)

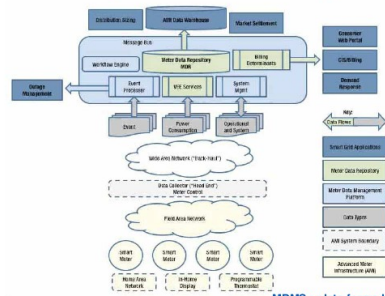
Fonctionnalités Smart-Grids : roadmap



Architectures pour les Smart-Grids

- ◆ L'expérimentation à grande échelle de projets Smart-Grids est une réalité (aux Etats-Unis)
 - Projets à forte composante informatique
 - Projets sociétaux (acceptabilité, privacy)
 - Exemples : PG&E (5M), SDG&E (1,5 M)
- ◆ Architectures mises en place :
 - Basées sur des solutions SI matures et standards (intégrations autour d'un bus d'entreprise)
 - Notion de MDMS (Metering Data Management System)
 - Fonctionnalités :
 - Canal batch pour les données horaires mises à disposition le lendemain sur portails web
 - Canal temps réel pour alertes
 - Mise à disposition des données du MDMS pour la facturation
 - Alimentation d'un entrepôt de données pour une historisation des données de comptage

Architectures pour les Smart-Grids (2)



MDMS : plate-forme logicielle qui acquiert des données de comptage depuis de nombreuses sources et les met à disposition, après intégration, synchronisation et nettoyage, auprès de nombreuses cibles.

Architectures pour les Smart-Grids (3)

- ◆ Architectures à moyen terme :
 - Elles intégreront une intelligence distribuée (routeurs intelligents à capacité de traitement, bases de données distribuées, traitements des données et événements temps réel (CEP) distribués) qui pourra être pilotée par un système centralisé
 - Traitement conjoint des données de comptage et des données liées au fonctionnement du réseau
 - Fonctions avancées de conduite et d'exploitation du réseau (observabilité, réseaux auto-cicatrisants ...)
 - Gestion de la demande, prix temps réels
 - Intégration étroite des énergies renouvelables
 - Optimisations locales
- **Systèmes de systèmes**

The necessary functions remain the same, the key issue is manage the complexity to support the necessary business capabilities at any scale as well as manage the separation of responsibilities to avoid "dueling" control systems (SDG&E)



Well Scalable Smart SigMA Adaptive of Massive Cloud

Travaux menés à EDF R&D

Données, infrastructures, technologies de traitement

Well Scalable Smart SigMA Adaptive of Massive Cloud



- ◆ Aujourd'hui :
 - Des panels de courbes de charges et des profils
 - Délai important de récupération et de mise à disposition des données
- ◆ Demain :
 - Des données de comptage à fine granularité, individuelles ou locales
 - Faible latence entre les systèmes opérationnels et les systèmes d'information décisionnels
- ◆ Des technologies de traitement de données « massives » permettant d'envisager une exploitation à grande échelle des données issues des infrastructures communicantes :
 - Gros entrepôts de données, décisionnel temps réel
 - Traitements de données à la volée (« Complex Event Processing », flux de données)
 - Traitements et décisions locales

Stockage Massif de Courbes de charge



◆ **Objectif** : montrer la **faisabilité** d'un stockage massif de courbes de charges rendues disponibles pour un certain nombre de traitements (plus ou moins complexes, plus ou moins concurrents, avec une latence variable selon les besoins)

- Données : courbes individuelles, données météo, informations contractuelles, données réseau
 - 1 relevé pas 10 mn / 35 millions de clients / an
 - Volumes de données annuels
 - 1800 milliards de relevés ; 600 To de données non compressées
 - Volumes de données journalier
 - 5 milliards de relevés ; ~2 To



- **Grandes fonctions évaluées**
 - Alimentation des données (par batch & par flux)
 - Rapprochement avec les données CRM
 - Pré-traitement des données
 - Synchronisation temporelle
 - Détection et correction d'anomalies
 - Changement de modes de représentation
 - Traitements
 - Calcul de synchrones par sous-population
 - Simulation de facturation (agile)
 - Simulation de traitements simultanés
 - Sélection d'une course sur motifs
 - Requêtes analytiques, BI
- **Critères d'évaluation**
 - Hautes performances
 - Alimentation des données, accès aux données, répartition des charges
 - Haute disponibilité
 - Tolérance aux pannes, redondance et secours
 - Extensibilité
 - Fonctionnelle et horizontale
 - Type d'architecture : share-nothing ou everything



Stockage Massif de Courbes de charge



◆ **Approches envisagées** :

- Approches distribuées :
 - POC interne, basé sur l'éco-système Hadoop en cours
 - Stockage et traitement distribué de type Map/Reduce
 - La maturité industrielle de ces approches n'est pas encore complètement atteinte
- Approches relationnelles classiques :
 - Travaux réalisés avec des partenaires
 - Appliances - POCs en cours
- Data Historian



◆ **Réalisation d'un générateur de données CourboGen**

- Permet la génération de courbes de charge et de données associées
- Outil paramétrable : pas de temps, durée, qualité des données, bruit sur les courbes
- Architecture distribuée
- Performances : 60000 relevés / seconde / module
- Sortie sous forme de flux de données



Constitution de courbes synchrones

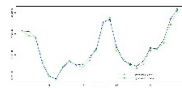


◆ **Objectif** : exploitation / valorisation des données de comptage massives pour différents besoins réseaux, avec une faible latence

- Constitution de courbes synchrones
 - Par segment logique (portefeuille client)
 - Par segment topologique (maillage réseau)
- ◆ **Points à résoudre**
 - Effectuer le calcul en temps réel, ou avec une grande réactivité
 - Volumes et bande passante
 - Qualité des données (pertes d'informations, retards, dysfonctionnements)

◆ **Éléments de réponse** :

- Architectures multi-chemins et mécanismes robustes à la duplication (= sketches ->)
- Architectures possibles et distribution de finalité ?
- Latence acceptée pour les traitements ?
- Inscrire ce calcul dans des outils de traitement de flux de données (traiter avant de stocker, à la volée)
 - Maquette utilisant le CEP StreamBase : traitement à la volée de 300000 courbes dupliquées 6 fois ; simulation de 2 RE (95% et 5% du volume des consommations)



Prévision de production locale



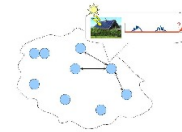
◆ **Objectif** : pilotage de l'équilibre du réseau au niveau local (intégration des ENR, gestion des VE), optimisations locales de fonctions réseau

◆ **Un élément de réponse** : le data-mining distribué

- Approches basées sur des agents indépendants les uns des autres qui peuvent effectuer des tâches différentes et collaborer ensemble
- Données très distribuées, vision globale sans centraliser l'ensemble des données
- Besoins de fouille locale, avec un contexte global (collaboration entre les tâches locales)

◆ **Premières expérimentations** :

- Prévision de production PV locale court terme



Prévision de consommation



◆ **Objectif** : prise en compte dans les méthodes et outils de prévision de l'évolution « smart grid » (données disponibles de façon détaillée : quelle opportunité ? Besoins de prévisions à mailles plus locales, et avec une plus grande variabilité des périmètres concernés)

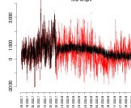
◆ **Aujourd'hui** :

- Les modèles de prévision sont robustes et de (très) bonne qualité
- Ils sont néanmoins conçus pour travailler sur un périmètre constant et un signal de consommation régulier

◆ **Demain** :

- Périmètres variables, signaux non stationnaires (mouvements de clients, périodes difficiles à prévoir, incertitudes, intégration de la production décentralisée, nouveaux usages)
- Quels bénéfices peut-on tirer de l'utilisation de données individuelles massives ?
- Données arrivant sous forme de flux, à prendre en compte pour les modèles court terme (de la veille pour le lendemain) et infra-journaliers
 - Adaptativité
 - Mouvements de clients captés avec faible latence
- **Données individuelles** :
 - Agrégation (clustering, échantillonnage)

Résultats de prévisions sur un an (fonction de taille à 88, modification du niveau moyen et de l'arrêt relatif en hiver)
RIVESE GAN : 908 MW, RIVESE GAN en Espagne : 714 MW



EDF R&D : Créer de la valeur et préparer l'avenir

Conclusion



Données massives pour les smart grids : impacts et enjeux

- ◆ Les projets « Smart-Grids » ont une forte composante informatique
 - Un réseau informatique, et pas seulement électrique
 - NTIC et utilities
 - Une augmentation forte du volume des données à traiter sur lesquelles on va vouloir appliquer des traitements parfois complexes
 - Place majeure du / des SI (volumétrie, complexité, architecture et urbanisme, agilité)
- ◆ Passage à l'échelle et verrous :
 - Stockage de larges volumes de séries temporelles
 - Approches centralisées et/ou distribuées
 - Traitement de données à la volée et en temps réel (CEP éventuellement distribués)
 - Data-mining à large échelle :
 - Distribution et parallélisation des algorithmes de fouille (cloud)
 - Le data-mining distribué peut amener des réponses en terme de volumétrie et de respect de la privacy
 - Apprentissage en ligne



Travaux réalisés avec :

Charles Bernard
 Alexis Bondu
 Xavier Brossat
 Youssa Chabchoub
 Leeley Daio Pires Dos Santos
 Alzennyr Gomes Da Silva
 Veronica Gomez
 Yannig Goude
 Benoît Grossin
 Georges Hébraïl
 Bruno Jacquin
 Sylvie Mallet
 Amandine Pierrot
 David Worms



Rererences

- ◆ **MDMS** : rapport du GTM Research
www.gtmresearch.com/rapport/the-emergence-of-meter-data-management-mdm
- ◆ **Smarter Energy @ IBM Research**, Brian Gaucher, Manager Smarter Energy, IBM T.J. Watson Research Center
- ◆ **Analytics and transactive control design for the Pacific Northwest Smart Grid Demonstration Project**, P. Huang, J. Kalagnanam, R. Natarajan (IBM Research Watson), D. Hammerstrom and R. Melton, Battle Memorial Institute, Pacific Northwest Division, Richland. (<http://www.ieee-smartgridcomm.org/techprogram.html>)
- ◆ hadoop.apache.org
- ◆ **Agrégation robuste de données en temps réel : application aux compteurs électriques communicants**, Y. Chabchoub, B. Grossin, Prix du meilleur article applicatif EGC 2011
- ◆ **Short term electricity load forecasting with adaptive GAM models**, Y. Goude, A. Pierrot, ISF 2010
- ◆ **A range of methods for electrical consumption forecasting**, X. Brossat, ISF 2010.
- ◆ « **SMAC – Stockage Massif de Courbes** », séminaire BILab du 12 Mai 2011, B. Jacquin, L. Daio Pires Dos Santos, A. Gomes Da Silva, D. Worms
- ◆ <http://www.csee.umbc.edu/~hillol/Kargupta/pubs.html>



2.10 David KONOPNICKI (IBM Haïfa-Research)

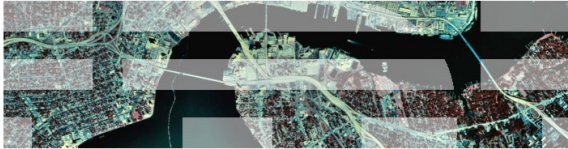
Massive-scale Analytics for a Smarter Planet

Everyday, we create 2.5 quintillion bytes of data—so much that 90% of the data in the world today has been created in the last two years alone. This data comes from everywhere : from sensors used to gather climate information, posts to social media sites, digital pictures and videos posted online, transaction records of online purchases, and from cell phone GPS signals to name a few. This data is Big Data. Big Data is more than a challenge ; it is an opportunity to find insight in new and emerging types of data and to answer questions that, in the past, were beyond reach. Until now, there was no practical way to harvest this opportunity. Today, IBM's platform for Big Data opens the door to a world of possibilities, giving organizations a solution that is designed specifically with the needs of the enterprise in mind and provides the infrastructure of a Smarter Planet : intelligence is being infused into the systems and processes that make the world work—into things no one would recognize as computers : cars, appliances, roadways, power grids, clothes, even natural systems such as agriculture and waterways.

David Konopnicki - Haifa Research Lab



Massive Scale Analytics for a Smarter Planet

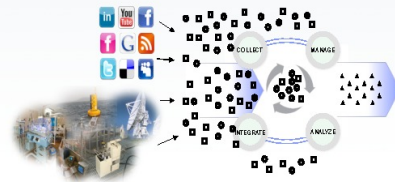


© 2009 IBM Corporation



The Big Data Challenge

- Manage and benefit from massive and growing amounts of data
 - 44x growth in coming decade from 800,000 petabytes to 35 zettabytes
- Handle unstructured data (text/images/video), social (graph) data...
- Exploit **BIG Data** in a timely and cost effective fashion

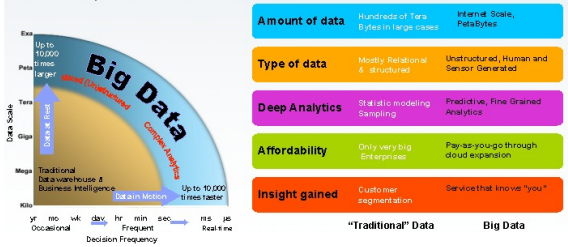


© 2009 IBM Corporation

Massive Scale Analytics



- A new set of tools is emerging around Big Data
- Internet technology for handling such Big Data is now entering the enterprise

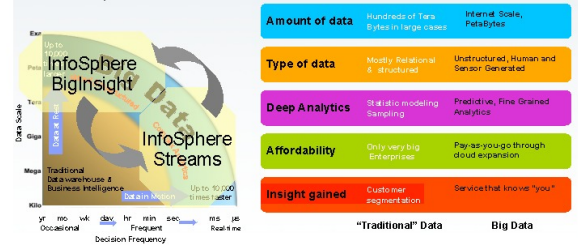


© 2009 IBM Corporation

Massive Scale Analytics



- A new set of tools is emerging around Big Data
- Internet technology for handling such Big Data is now entering the enterprise



© 2009 IBM Corporation

Massive Scale Data Analytics Problem and Solution



- How do you **scale up** applications?
 - Run jobs processing 100's of terabytes of data
 - Takes 11 days to read on 1 computer
- Need lots of cheap computers
 - Fixes speed problem (15 minutes on 1000 computers), but...
 - Reliability problems
 - In large clusters, computers fail every day
 - Cluster size is not fixed
- Need common infrastructure
 - Must be efficient and reliable
- Infrastructure:**
 - Google File System – *OSDP/2003*
 - Hadoop File System**
- Computing Paradigm:**
 - Map-Reduce** – *OSDP/2004*
 - Hadoop Map-Reduce**
- Scripting Language:**
 - Sawzall** – *Scientific Programming Journal/2005*
 - Pig, Jaql**
- NoSQL DB:**
 - Big Table** – *OSDP/2006*
 - HBase, Cassandra, Hive**
- Workflow:**
 - Oozie, Metatracker**

© 2009 IBM Corporation

Solution

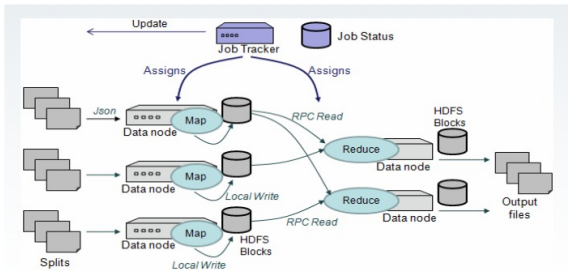


- Open Source Apache Project
- Hadoop Core includes:
 - Distributed File System - distributes data
 - Map/Reduce - distributes application
- Written in Java
- Runs on
 - Linux, Mac OSX, Windows, and Solaris
 - Commodity hardware

© 2009 IBM Corporation



Hadoop in a Nutshell



•An open-source computing platform that is both **distributed and redundant**, handles **structured and unstructured data** and supports a **simple and efficient programming paradigm** called map-reduce.

© 2009 IBM Corporation

© 2009 IBM Corporation

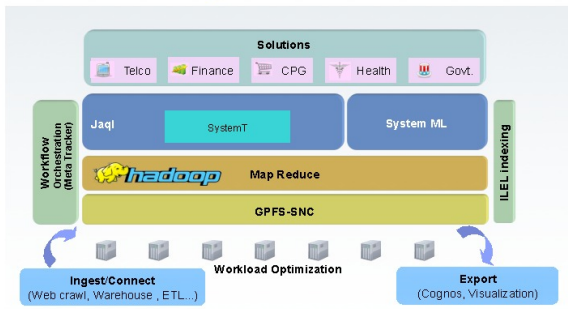


Who Uses Hadoop?

- Amazon/A9
- AOL
- Facebook
- Fox interactive media
- Google / IBM
- Netflix
- New York Times
- PowerSet (now Microsoft)
- Quantcast
- Rackspace/Mailtrust
- Veoh
- Yahoo!
- More at <http://wiki.apache.org/hadoop/PoweredBy>



IBM's Massive Scale Analytics Architecture



Innovation above, underneath and around Hadoop

© 2009 IBM Corporation



Products Impact

- Cognos Consumer Insight: Mining social media helps marketing professionals transform customer relationships by analyzing sentiment, affinity and evolving topics in social media sites.
- Infosphere BigInsight: Bringing the power of Hadoop to the Enterprise



© 2009 IBM Corporation



Smarter Planet: Something profound is happening...



INSTRUMENTED

We now have the ability to measure, sense and see the exact condition of practically everything.



INTERCONNECTED

People, systems and objects can communicate and interact with each other in entirely new ways.



INTELLIGENT

We can respond to changes quickly and accurately, and get better results by predicting and optimizing for future events.



© 2009 IBM Corporation



Intelligent systems that gather, synthesize and apply information will change the way entire industries operate.

Smart water

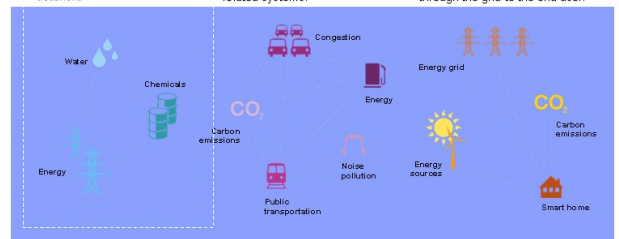
Apply monitoring and management technologies to help optimize the availability, delivery, use, and quality of water as well as related systems including energy and chemical treatment.

Smart traffic

Use real-time traffic prediction and dynamic tolling to reduce congestion and its byproducts while positively influencing related systems.

Smart energy

Analyze customer usage and provide customized products and services that help to boost efficiency from the source through the grid to the end user.



© 2009 IBM Corporation

How much water do you need to..

- 20 liters of water to manufacture one sheet of paper
- 40 liters of water to create a loaf of bread
- 70 liters of water to create one apple
- 80 liters of water per dollar of an industrial product
- 91 liters of water to manufacture one plastic can
- 120 liters of water to create a glass of wine
- 140 liters of water to manufacture one cup of coffee
- 1300 liters of water to manufacture one kilo of wheat
- 15,500 liters of water to create one kilo of beef
- 10855 liter of waters to manufacture one pair of jeans

IBM SMARTER PLANET

The world suffers from water shortage

"Water is the 21st century's oil"
Business Week

- Today, 1 out of 5 people does not have access to fresh water.
- By 2020, 1.3 billion people will not have access to fresh water.
- 50% of the world's population is projected to live in areas of water scarcity by 2030.
- By the year 2050, when the world's population will be 9.4 billion, water will become the most scarce resource.

Even Though, there's enough water for everyone: about 2 trillion liters of fresh water per person, when each consumes an average of 3 liters per day.

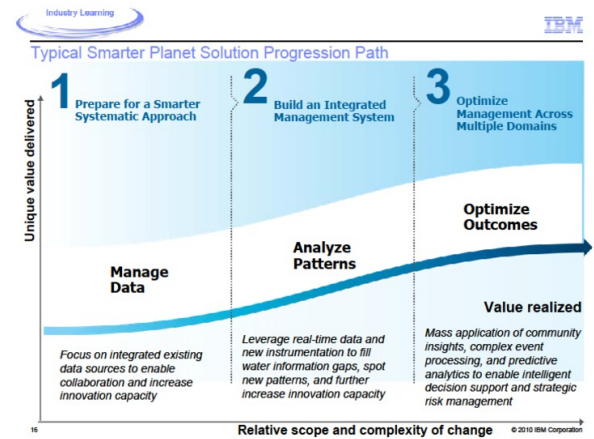
IBM SMARTER PLANET

Water shortage is an outcome of inefficiency and waste

Industry	Agriculture	Private consuming
22% of the world's water consumption is industrial use. In the last 100 years, water consumption has increased twice as much as the population increase.	Worldwide agriculture wastes about 60% of the 2500 trillion liters it consumes annually.	About 50% of water consumption is lost due to leakage and corrosion management.

The problem is not the lack of water, but managing water information about it's usage, that causes inefficiency and great waste.

IBM SMARTER PLANET



Addressing Non-Revenue Water using Analytics and Optimization

Leakage or Theft Detection at the Residential, public buildings.

Understand usage patterns and detect anomalies for low and high consumption to detect leakage, theft or faulty meters

Leakage Detection at the Network Level using optimization

Find "optimal" location of leak(s) to explain difference between actual measurements and model predicted measurements

Leakage Reduction using Dynamic Pressure Control

Create optimization model to adjust the pressure dynamically so that only the required flow will be supplied yielding cost reduction in energy and water achieved.

Optimal Valve Placement for Pressure Reduction

Use an optimization model to find the optimal number of valves, and their location, so as to enable the most effective pressure management

Why now? "Because we must"

Leaders focus on credit crisis
Roubini warns of double-dip recession: rep...
Consumers Postpone Purchases as Recession Deepens
One certainty: Gas will go up
Where will it end?
World Bank warns of social unrest
Why The Telcos Are Doomed
WHAT WENT WRONG WITH ECONOMICS?
Climate Change When glaciers start moving
U.S. must "adapt" to changing climate
Can Digital Health Protect Your Privacy?
Enemies in the media

IBM SMARTER PLANET

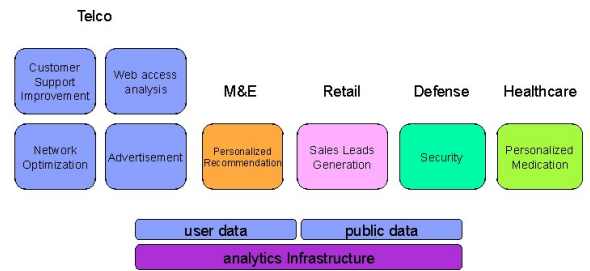


IBM Customer Modeler

IBM Customer Modeler

- A components library to analyze customer behavioral data and enable new insights and business scenarios based on web behavioral data
- Behavioral data:
 - Web Search activities (internal and external)
 - Web Browsing (internal and external)
 - Customer generated content (blogs, twitter)
 - Products browsed, bought
 - Location, context
- Other sources:
 - Enterprise:
 - Product Catalogs
 - Web site metadata
 - Customer data, segmentations, taxonomies
 - Public:
 - Wikipedia
 - Open directory projects

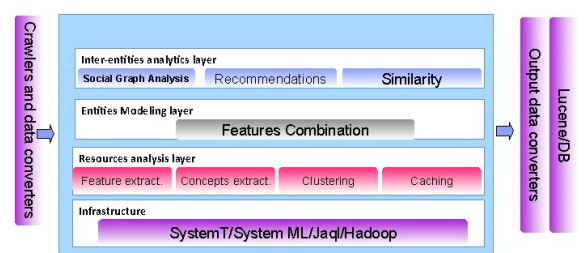
IBM Customer Modeler in Industries



Scenarios

- Network Optimization: Learning customer browsing activities in order to predict future data access and optimize caching
- Advertisement: Learning customer interests through browsing activities in order to optimize ads targeting
- Customer support improvement: Learning customer search and browsing activities in order to discover information need and improve interaction, knowledge bases
- Web access analysis: Learning customer interests and clustering to detect trends
- Personalized Recommendation: Learning customer preferences from consumption data and providing personalized recommendations
- Sales leads generation: Analyzing customers (private and companies) generated content to discover opportunities
- Healthcare: Analyzing social media to find data relevant to particular patients

Biginsight Library for Industry Scenarios

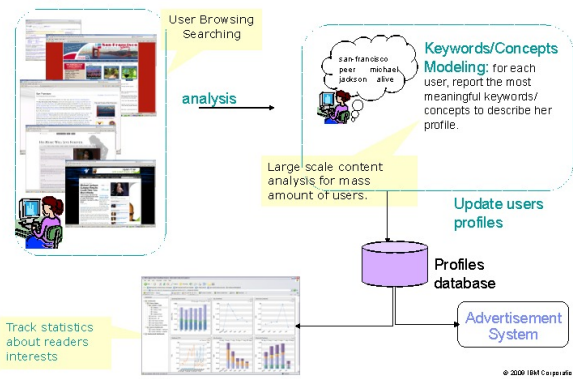


Components in Details

Components

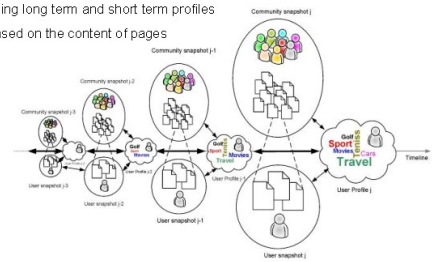
- Web Browsing/Searching Analytics
- Personalized Recommendations
- Finding Influencers
- Finding Similarities

Mining Web Behavior

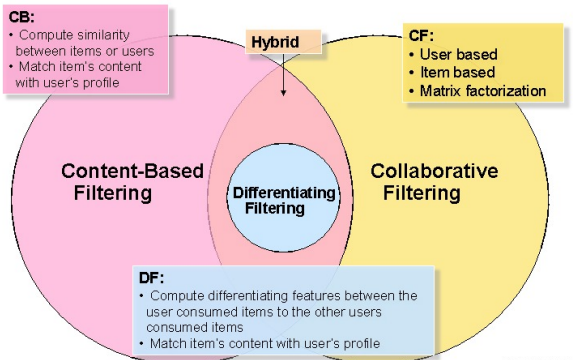


Mining Web Behavior

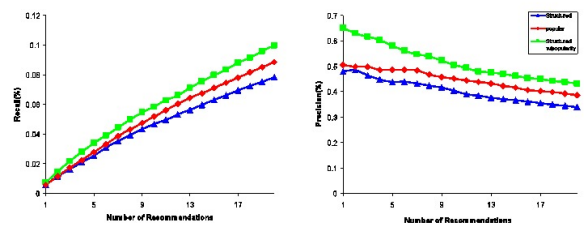
- Source: Internal/External Browsing/Searching History
- Basic concept: finding what statistically differentiates a particular user from the rest of the population
- Maintaining and combining long term and short term profiles
- Can be implemented based on the content of pages or based on metadata



Preferences Prediction for recommendations



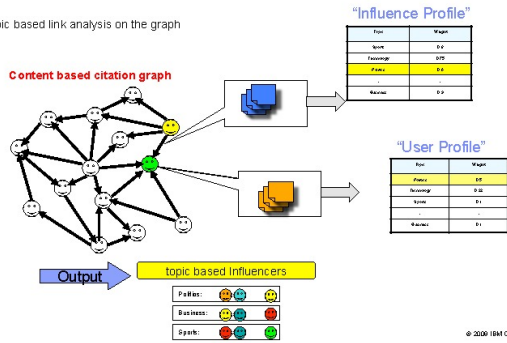
Enhanced Popularity



> Precision over popularity improved by 28% at 1, 19% at 5

Finding Influencers in Social Media

- Build user profiles and "link" profiles (Blogs, Twitter...)
- Topic based link analysis on the graph

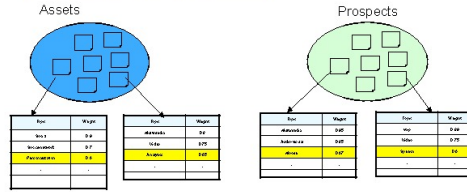


Similarity

Example: rank prospects by decreasing similarity to assets

- > Use assets and prospects profiles
- > Apply cos-similarity on each (a,p) pair
- > Take for each asset the top-k prospects

$$score(a, p) = \sum_{t \in Profile(a) \cap Profile(p)} w(t, a) * w(t, p) * idf(t)$$



Summary

- Using massive scale analytics to improve IBM customers business scenarios
 - analyzing more data (e.g., online behavior)
 - analyzing new types of data (e.g., text, social relationships)
 - developing finer grained models (personal models)

<http://www.association-aristote.fr> info@association-aristote.fr

ARISTOTE Association Loi de 1901. Siège social : CEA-DSI CEN Saclay Bât. 474, 91191 Gif-sur-Yvette Cedex.
Secrétariat : Aristote, École Polytechnique, 91128 Palaiseau Cedex.
Tél. : +33(0)1 69 33 99 66 Fax : +33(0)1 69 33 99 67 Courriel : Marie.Tetard@polytechnique.edu
Site internet <http://www.association-aristote.fr>