# Le management d'une immense banque de données

## Stéphane Mouton

*SST Department manager*

16ème rencontre EVELINE M ARKIEWICZ

Brussels, 14/11/2015

Centre d'Excellence en **Technologies** de
l'**Information** et de la **Communication**

www.cetic.be

# In the beginning …

# Then IT resource became affordable

16ème rencontre Eveline Markiewicz

# But data production explodes

# What is the status now?

- Network lags behind

"Never underestimate the bandwidth of a station wagon full of tapes hurtling down the highway."

Tanenbaum, Andrew S. (1989). Computer Networks

Fiber network is crucial

# What is the status now?

- Most of "usual" commercial offering don't stack up seamlessly

16ème rencontre Eveline Markiewicz

# Need for scalability

- The capability of IT systems to adapt to system demand at a predictable, ideally linear, pace.

# Scalable storage to cope with volume

- When a single IT resource can't handle all the data anymore
  - *Distributed* solutions, based on many cloned (and possibly inexpensive) *nodes*



10GbE public access network
10GbE front-side storage network
10GbE back-side storage network
1GbE management network

# Scalable databases

- New ways to structure data storage
    - NoSQL (Not Only SQL databases)



**BDD Orientée colonnes**

**BDD Clé-Valeur**

Relational
Key-Value
Column-Oriented
Document-Oriented

**BDD Orientée document**

**BDD Orientée graphe**

Consistency — RDBMs (Oracle, MySQL), Aster Data, Green Plum, Vertica — Availability

Pick Any Two

mongoDB, Terrastore, Datastore, Hypertable, Hbase, Redis, Berkeley DB, MemcacheDB, Scalaris

Dynamo, Voldemort, Tokyo Cabinet, KAI, Cassandra, SimpleDB, CouchDB, Riak

Partition tolerance

# Process and analyse data in a scalable way

- What if you have possibly an infinite number of computers to process your data?
  - Your data analysis doesn't change, the way you express it does.

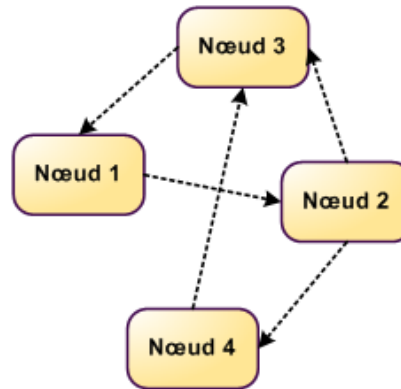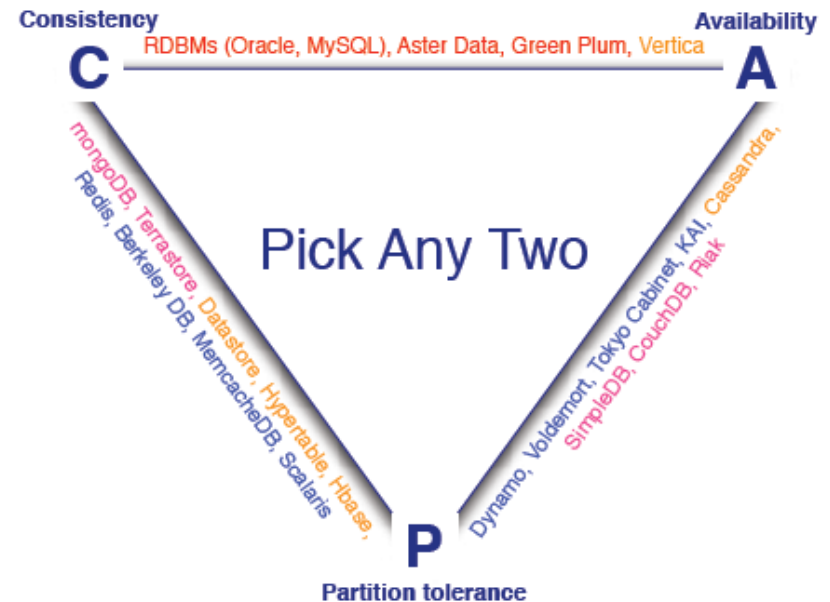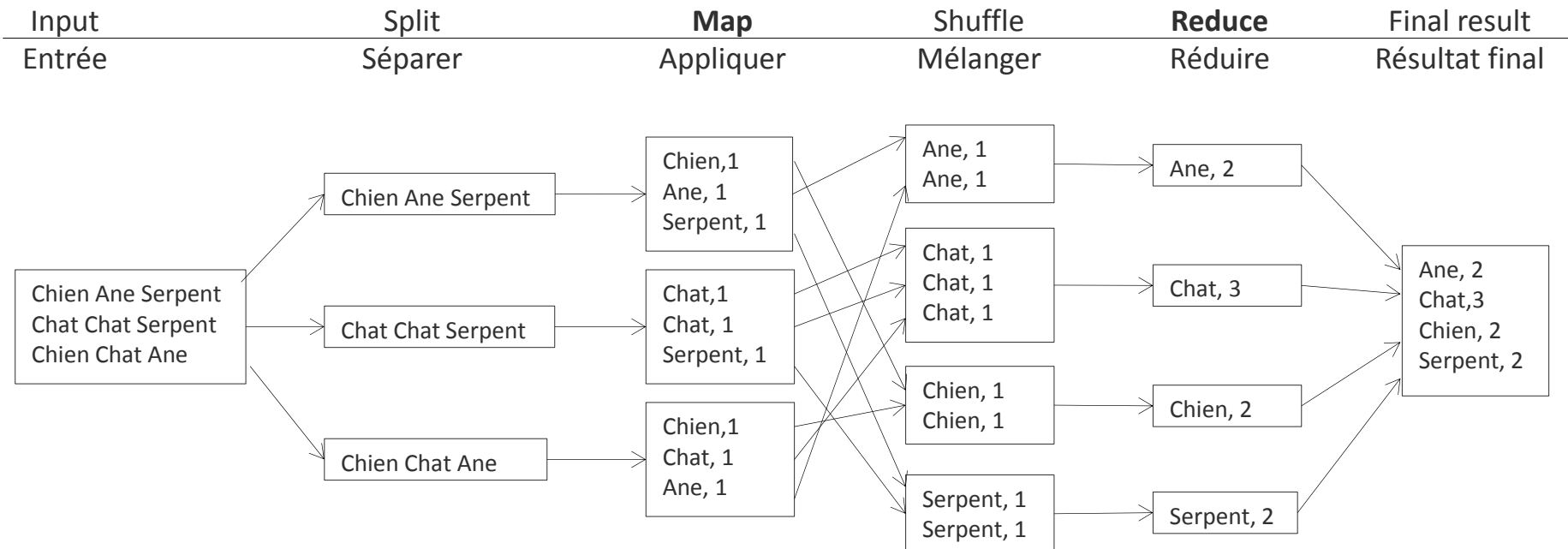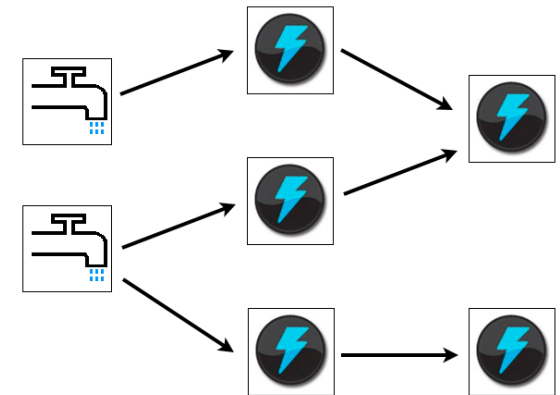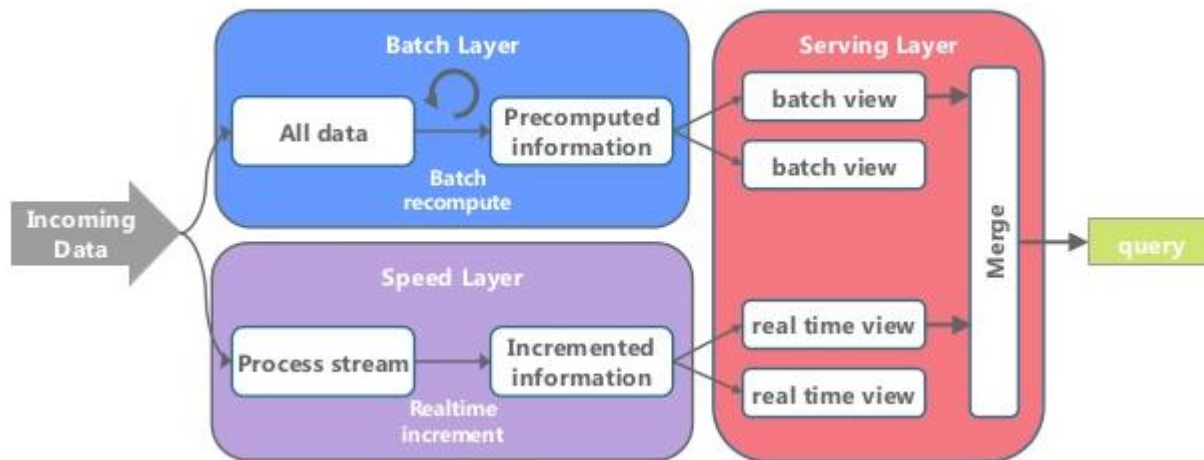| Input | Split | **Map** | Shuffle | **Reduce** | Final result |
|---|---|---|---|---|---|
| Entrée | Séparer | Appliquer | Mélanger | Réduire | Résultat final |

Chien Ane Serpent Chat Chat Serpent Chien Chat Ane

Chien Ane Serpent

Chat Chat Serpent

Chien Chat Ane

Chien,1
Ane, 1
Serpent, 1

Chat,1
Chat, 1
Serpent, 1

Chien,1
Chat, 1
Ane, 1

Ane, 1
Ane, 1

Chat, 1
Chat, 1
Chat, 1

Chien, 1
Chien, 1

Serpent, 1
Serpent, 1

Ane, 2

Chat, 3

Chien, 2

Serpent, 2
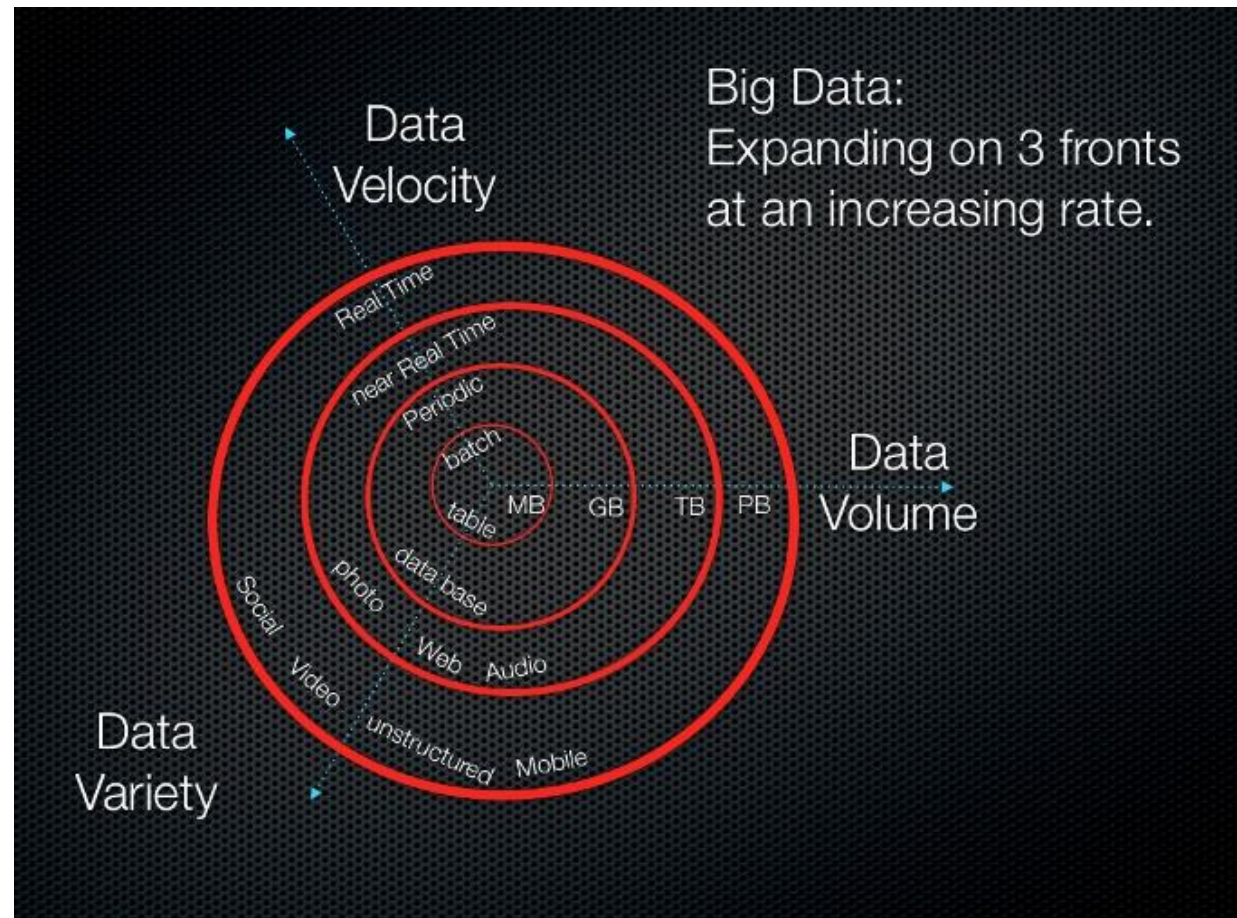
Ane, 2
Chat,3
Chien, 2
Serpent, 2

# Process and analyse data in a scalable way

- New programming frameworks for data analysis

# Big Data

- A relative concept, linked to data
  - Volume
  - Variety
  - Velocity

# Volume: data integration approaches

- 3 main models ensuring data bases autonomy
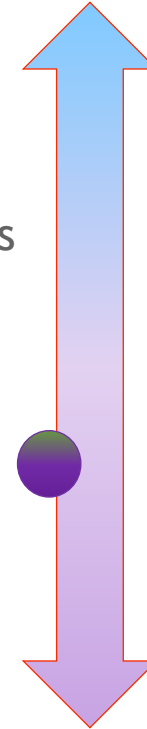  - Full Data Integration (Data Warehousing)
    - Not suitable for data with a very short-term relevance
    - Can lead to improperly formed data warehouses, not focused on the analysis topic

  - Federated databases
    - One global schema
    - One main data format/model
    - Mostly static

  - Mediator/Wrapper
    - Support dynamic sets of data sources
    - High flexibility: schema can evolve at query time
    - Wrappers to handle heterogeneous data format

**Tightly coupled**
- Administrator builds a unique FDB schema

**Loosely coupled**
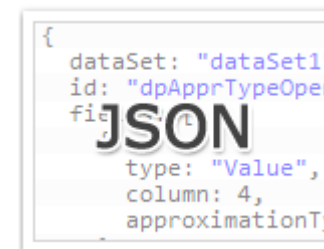- Component DBs out of control
- Import schemas as view over data sources

# Data integration leads to data variety
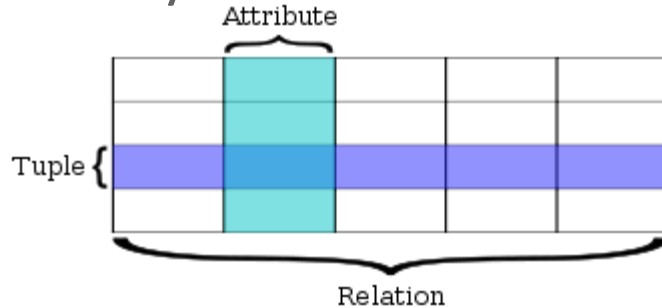
- Variety in content type

16ème rencontre Eveline Markiewicz

# Data integration leads to data variety

- Variety in data format

16ème rencontre Eveline Markiewicz

# Data integration leads to data variety

- Variety in data model



```xml
<Books>
    <Book ISBN="0553212419">
        <title>Sherlock Holmes: Complete Novels...
        <author>Sir Arthur Conan Doyle</author>
    </Book>
    <Book ISBN="0743273567">
        <title>The Great Gatsby</title>
        <author>F. Scott Fitzgerald</author>
    </Book>
    <Book ISBN="0684826976">
        <title>Undaunted Courage</title>
        <author>Stephen E. Ambrose</author>
    </Book>
    <Book ISBN="0743203178">
        <title>Nothing Like It In the World</title>
        <author>Stephen E. Ambrose</author>
    </Book>
</Books>
```

# Data integration leads to data variety

- Variety in data source

# Challenges of data variety

- Identify relevant data sources

- Extract, clean and store data

- Understand the data (semantics, meta-data)

- Deliver information

# Data management also applies

- Not only in terms of cost

# Questions?

**cetic**
Your Connection to **ICT** Research

Aéropole de Charleroi-Gosselies
*Bâtiment Éole*
Rue des Frères Wright, 29/3
B-6041 Charleroi

Tel: +32.71.490.700
Fax: +32.71.490.799

**www.cetic.be**
**info@cetic.be**

linkedin.com/company/cetic

twitter.com/@CETIC

**Stéphane Mouton**
*SST Department manager*

Tel : +32 71 490 726

Mob : +32 475 76 78 50

stephane.mouton@cetic.be