# THALES

# Technological Challenges for Multimodal Information Mastering
## *Workshop Franco-Finlandais 5/5/06*

*François Marcotorchino*

Systèmes Terre et Interarmées

- **Data Growth**
  - **Quantity of available data worldwide doubles every year**
- **Availability of processing tools**
  - **Cost of data storage has decreased from €/Mb to cents / Mb**
  - **Computing power continues to double every year**
- **3 Ages of the Information Era**

**2000 and beyond**

*Understand*

**Knowledge Engineering**

**90s and beyond**

*Carry…*

**Telecommunications**
- **Bandwidth**
- **Mobility**
- **Security**

**80s and beyond**

*Process…*

**Computing**
- **Computing power**
- **Data storage**

**THALES**

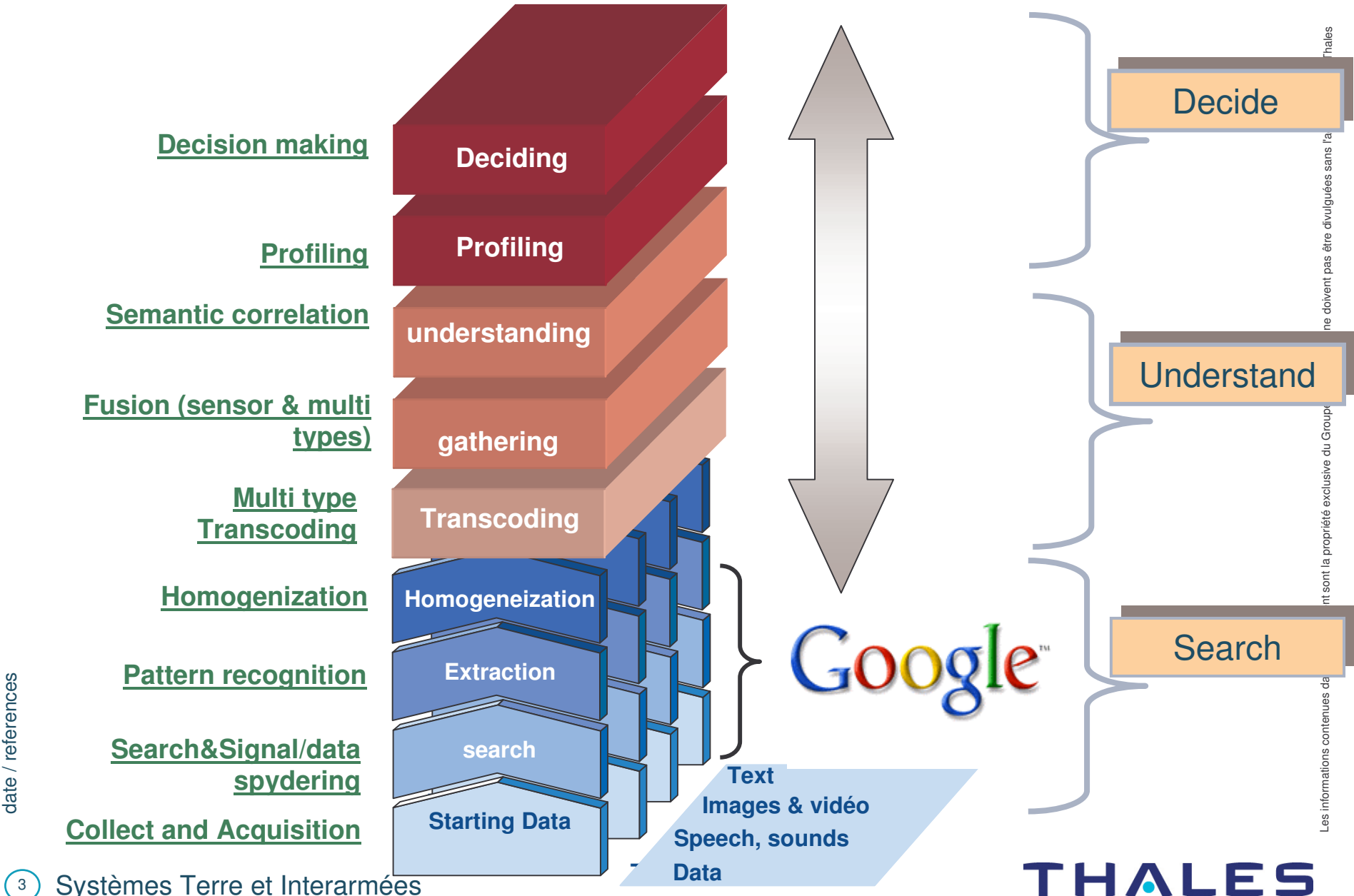# Technological   Challenges of Multimodal Information Mastering

- **Be at the level of the best for the indexation of the  web**

- **Improve the current processes of consultation**

- **Drastically improve "Miners*" technology**

- **Merge and inter-correlate multi-source data**

- **Achieve breakthrough in « machine automata » technology to improve *extraction* quality**

- **Optimize the transcoding process**

- **Achieve true interoperability between technologies**

- **Manage multimedia bases with the same efficiency as current text or digital data databases**

- **Develop applications in the Domains of : Security, Health, E-training, Digital libraries, Finance, and  Digital Life**

**THALES**

**Decision making** — Deciding

**Profiling** — Profiling

**Semantic correlation** — understanding

**Fusion (sensor & multi types)** — gathering

**Multi type Transcoding** — Transcoding

**Homogenization** — Homogeneization

**Pattern recognition** — Extraction

**Search&Signal/data spydering** — search

**Collect and Acquisition** — Starting Data

Text
Images & vidéo
Speech, sounds
Data

Google

Decide

Understand

Search

date / references

THALES

- **Information Extraction**

- **Automatic  Clustering (Structured/unstructured)**

- **Transcoding Data from Type A to Type B**

- **Information Fusion**

- **Semantic Graphs**

**THALES**

## Different Types of Data Fusion Type A➜Type B

■**Fusion: Text ➜Numerical Structured Data**

This fusion (or joint treatment Textual Data with Structured Data ) must be understood as semantic reciprocal treatments, allowing semantic correlation and not "simple juxtaposition"

■**Fusion : Speech ➜ Text (already under process)**

Transformation by means of transcoders of the "speech signal" into its "textual translation" (example : Voice Dictation IBM, LIMSI tools etc..)

■**Fusion : Image➜ Structured Numerical Data (more classical)**

Here we must use tools for "image contours and outlines analysis" for translate Image into "Semantic Descriptive Qualitative Markers"
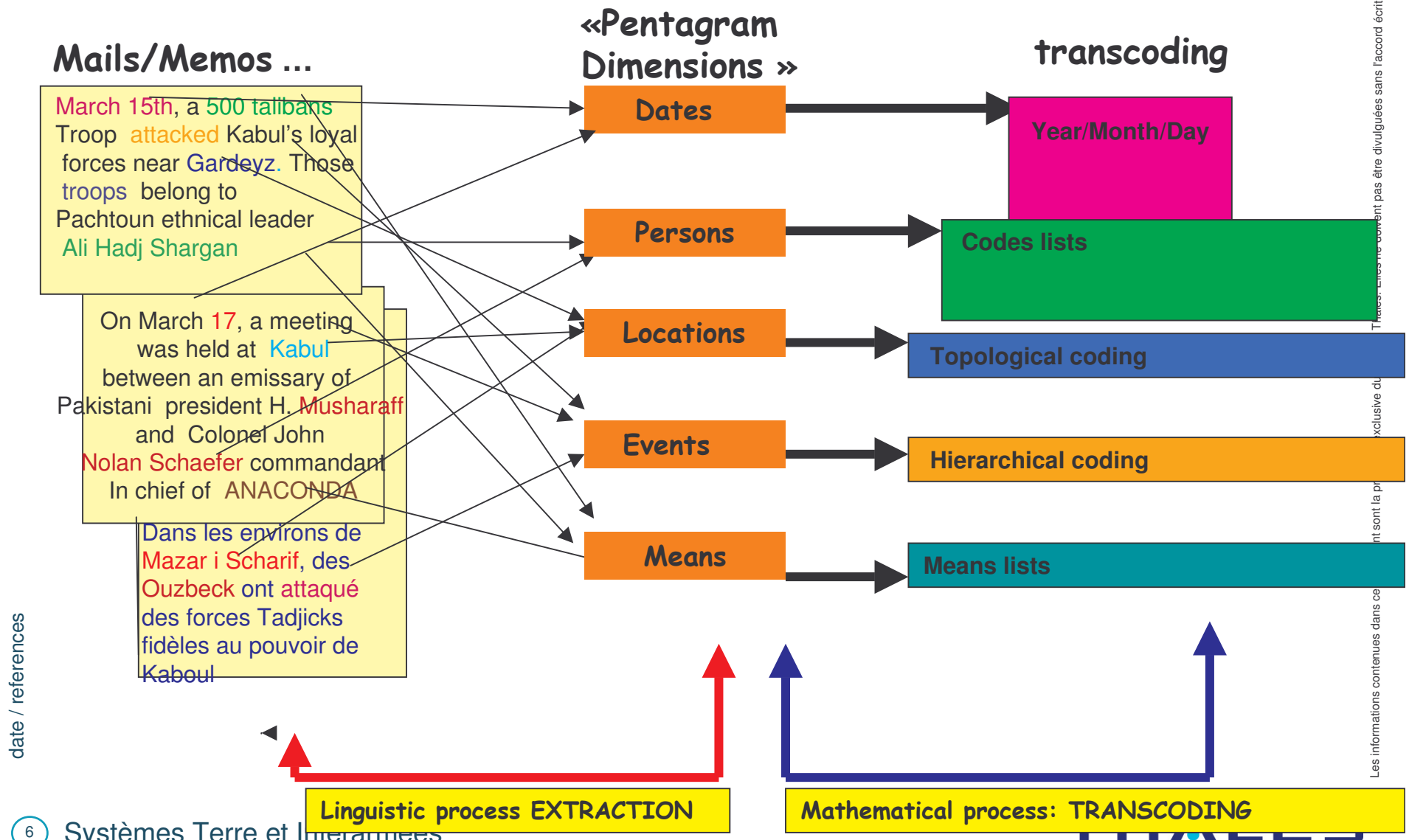
■**Fusion : Image ➜ Text**

Here we must try to interpret (if possible) the contours and forms by means of predetermined structures (gabarits) with associated generation of textual scripts.

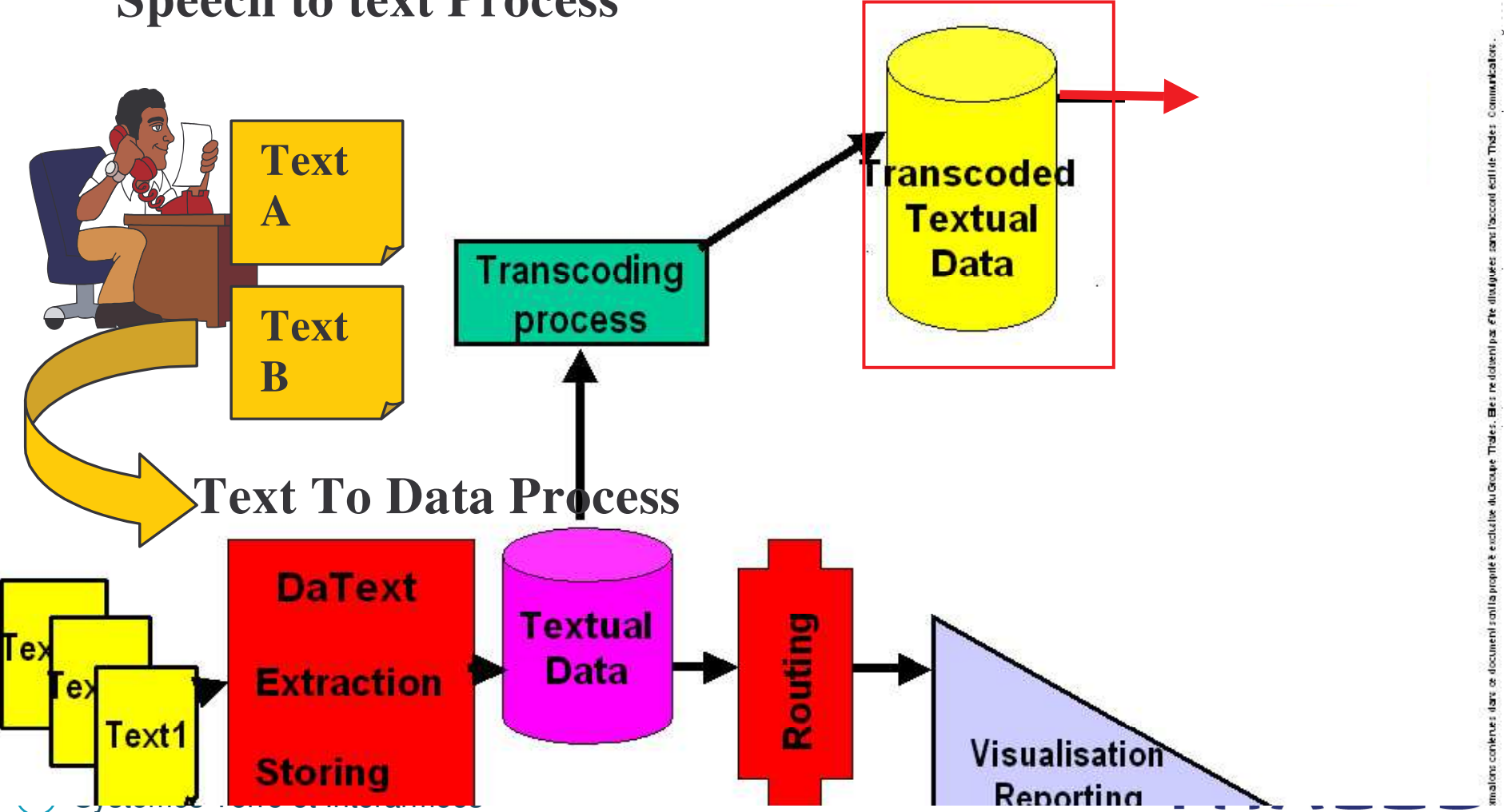■**Fusion: Cartographical Data ➜ Structured Numerical Data**

Dans In that case we translate the topographical Information into a set of Qualitative (hierarchical or Nominal) compatible with the other descriptive numerical variables already stored.

date / references

**THALES**

# Extracted information transcoding process:
## Text ➜ Data

**Mails/Memos ...**

**«Pentagram Dimensions »**

**transcoding**

March 15th, a 500 talibans Troop attacked Kabul's loyal forces near Gardeyz. Those troops belong to Pachtoun ethnical leader Ali Hadj Shargan

On March 17, a meeting was held at Kabul between an emissary of Pakistani president H. Musharaff and Colonel John Nolan Schaefer commandant In chief of ANACONDA

Dans les environs de Mazar i Scharif, des Ouzbeck ont attaqué des forces Tadjicks fidèles au pouvoir de Kaboul

**Dates** ➜ **Year/Month/Day**

**Persons** ➜ **Codes lists**

**Locations** ➜ **Topological coding**

**Events** ➜ **Hierarchical coding**

**Means** ➜ **Means lists**

**Linguistic process EXTRACTION**

**Mathematical process: TRANSCODING**

# Speech to text and Information Extraction Chain Coupling

Speech to text Process

Text A

Text B

Transcoding process

Transcoded Textual Data

Text To Data Process

Text1

DaText Extraction Storing

Textual Data

Routing

Visualisation Reporting

Exogen data (Pricings, Costs, sizes)

Simulation on model building

Structured Data Base X

Structured Data Base Y

Transcoded Textual Data

FUSION + Optimized Pre-Processing

Data & Text Analysis and Mining
Clustering
Segmentation
Seriation
Categorisation
Scoring

Visualisations Graphics Cartographies

Optimisation Processing Power Pricing Penalty Processing Logistic Opt. Etc..

Simulations

Simulation on data sets

Resulting Scenarios ROI Measures

THALES COM 12 02 2004

Thales Communications

THALES

Thales Communications.

Les informations contenues dans ce document sont la propriété exclusive du Groupe

**Numerical Data Base**

**Speech Data Base**

**Images Data Base**

**Textual** data

Direct without recoding

**Transforming the pitch and intensity of signal into new variables**

**Transforming the Image Outlines and contours into new variables**

**Through extraction from text of numerical Dimensions of discourse (Transcoding)**

**Merged Numerical Data Base**

**Exhaustive Treatment of the Full Data Base through powerful Data Mining Engines (up to 20 000 000 records in one shot)**

**First Strategy: Transform all the set of data into structured ones:**

**Very powerful for huge sized treatments but loss of semantical value**

**THALES**

date / references

**Numerical Data Base**

**Speech Data Base**

**Images Data Base**

**Textual data**

**Through "Speech to Text" Process**

**Through "Image to text" Process**

**Direct without recoding**

**Transcoding by automatic script generation or "Ideliance like" pre-structuring**

**Textual data**

**Exhaustive Treatment of the Full Text Repository through powerful Text Mining Engines (up to 1 000 000 records in one shot)**

**Second Strategy: Transform the whole data set into textual information : Very powerful for keeping a high level of semantic value, less powerful for large sized problems**

THALES

**From a SEMANTIC GRAPH, a large variety of Problems related to Clustering and Similarity is appearing:**

**A each node of the Graph one gives the attributes which are functions of the linked sub-graphs**

**Consider those sub-graphs per say -as new objects to classify (or to cluster)**

**Research all the chains of the graphs which constitute the attributes for their classification or clustering**

**The Clustering results are a real enrichment for the Semantic Graph and go directly into the TURBO machine for other treatments.**



STEP A: MERGE MANY SOURCES INTO AN UNIQUE SEMANTIC NETWORK OF MEETINGS, PEOPLE, PLACES, DATES ...



STEP B: AUTOMATIC DISCOVERY OF CLUSTERS OF OBJECTS SHARING PROPERTIES

What will Google do next? How will Yahoo counter the move? And how will this impact MSN Search, Ask Jeeves and the scores of smaller players vying to get their share of the search engine revenue pie?

## Clustering Catches On

« *Another common prediction is for* **the increased adoption of clustering technology**. *America Online is already offering clustering via its Vivisimo partnership.* **Gartner analyst Alan Weiner**, *for one, said he's hooked on clustering, and he* **expects to see major search engine players add clustering features in 2005** *to make search more user-friendly.*

*Clustering is an elemental way of taking people through a more direct path to what they are looking for," Gartner analyst said. "If you type in the world 'polish,' the search engine might not know if you are looking for information about Poland or products that make your car shiny. With clustering technology, you have on-the-fly categories and you can immediately choose 'car-related accessories. »*

**THALES**

May 23-24, 2006
Hilton New York, New York City

**Beyond Search: Intelligent Use of Intelligence**

*William Lunceford, Section Manager, Procter & Gamble*

« Search tools are more abundant than ever, but the people you support are spending valuable time sifting through irrelevant search results and may be missing critical information. How can you improve their overall search experience, and, more importantly, how can you help them make intelligent use of search intelligence? P&G created a unified search tool that sorts results **into clusters** that are intelligently selected from words and phrases found in the documents themselves. Learn about the benefits of clustering and how P&G's project evolved to extend across the entire enterprise ». *Hosted by Vivisimo.*

date / references

**THALES**

# « Sprint Platform™ » functionnal spectrum

- **Automatic treatments of Documents**

  - Indexation/categorization/classification plans (**pentagramm**)

  - Multi-modal Filtering, highlighting of « named entities »

  - Hierarchy of Documents according to « **gravity or criticality** » levels

  - Automatic Summaries

  - Reporting

- **Dynamic Cartography**

  - Automatic Geolocation of events

  - Geographical **Inverse search** of documents

  - Filtering and annotations

- **Decision Aided process**

  - **Clustering**

  - Semantic Networks

INFOM@GIC

THALES

date / references

INFOM@GIC Multimédia & Vie Numérique

■ *Knowledge Engineering* : **Develop worldclass**

**technologies enabling « Ile de France » companies to gain**

**lead positions on a fast growing market (> 20B€ 2008)**

■ **Capitalize on the unique scientific potential of Ile-de-**

**France companies, labs, universities, start-ups by**

**federating efforts**

**Enable everyone, private or professional, to *navigate* easily and smartly in huge amount of multi-source information so as to *understand* and *decide***
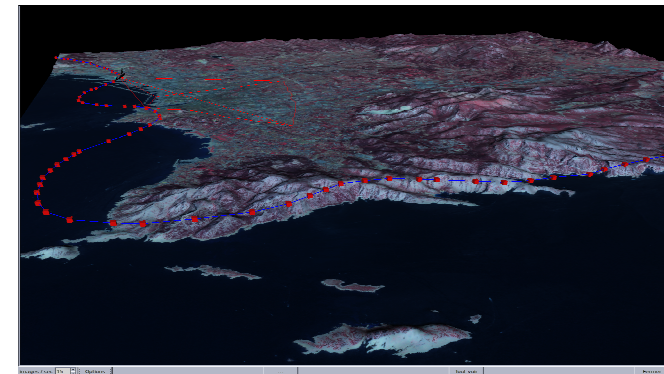
INFOM@GIC Multimédia & Vie Numérique

- **Etre au niveau** des meilleurs pour l'indexation du Web

- **Améliorer** les processus de **consultation** actuels

- Avoir de meilleurs « **Miners\*** » qu'IBM ou ORACLE

- **Fusionner et inter-corréler** les données de nature et sources différentes au delà de l'Etat de l'Art actuel

- **Dépasser** les limites actuelles des technologies d'automates pour gagner de la performance au niveau **Extraction**

- **Optimiser** les transcodages « **Type de données X vers Type de données Y** »

- Passer d'une simple juxtaposition à une réelle **interopérabilité** des outils.

- **Gérer des bases multimédias** avec la même efficacité que la gestion actuelle des bases textuelles ou de données numériques

- Trouver des applications **différenciatrices** dans les domaines de la Sécurité, de la Santé, de la Finance, et de la Vie Numérique

date / references

PARIS ILE-DE-FRANCE

# Infom@gic facts



- **<u>Budget</u> : 65M€**

- **<u>Duration</u> : 3 years**

- **<u>33 partners</u>**

- **<u>Calendar</u> :**

  - **2006 : models and implementation of first prototypes**

  - **2007 : complete suite of prototypes**

  - **2008 : technological integration platform and software bricks**



**<u>Budget split</u>**

- **Big firms and large Institutes : 58,6%**

- **Small and Medium Sized Companies : 21,2%**

- **Research Centres and Universities :20,2%**

date / references