SMR.1589 - 4

# Workshop on
# Managing Nuclear Knowledge

# 8 - 12 November 2004

----------------------------------------------------------------------------------------------------------------

# Searching and Accessing Information

**Heinz BACHMANN**
**CONVERA AG**
**Flawilerstrasse 27**
**9500 Wil**
**SWITZERLAND**

----------------------------------------------------------------------------------------------------------------

# CONVERA OVERVIEW

IAEA

Trieste - November 2004

# Agenda

- About CONVERA

- Introduction

- CONVERA Platform Infrastructure

- CONVERA's Retrieval Technologies

- Visualisation

- Automatic Categorization Technologies

- Dynamic Classification Technologies

- What we do / What's Third Party
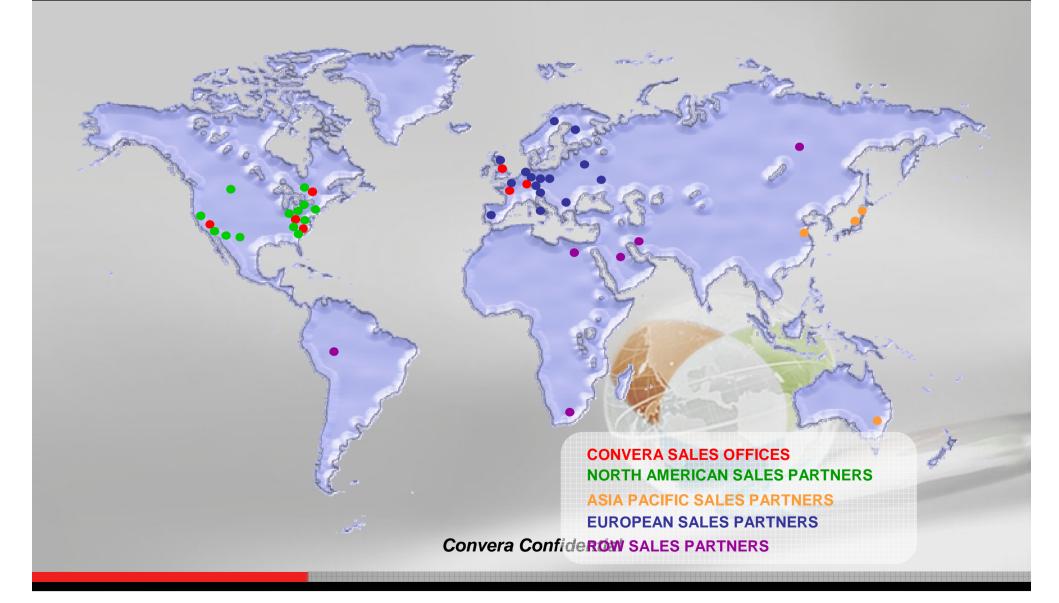
- Product demonstration

## Convera Corporation



**Convera is a leading provider of enterprise wide index, search and categorization software products and solutions**

- 20+ years of innovation in intelligent information infrastructure
- 250 employees
- 800 customers in 29 countries
- 70+ business partners
- Publicly traded (NASDAQ: CNVR)

*Convera Confidential*

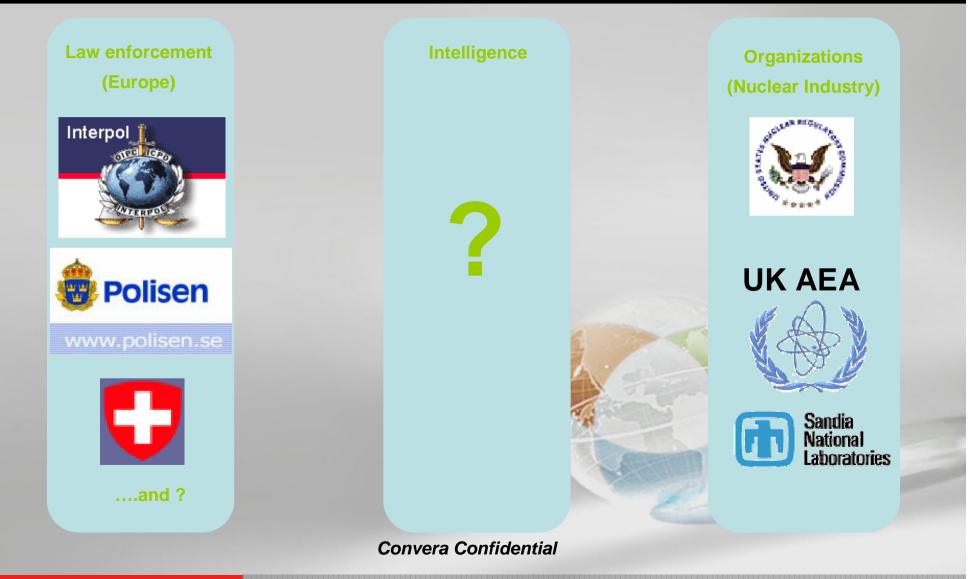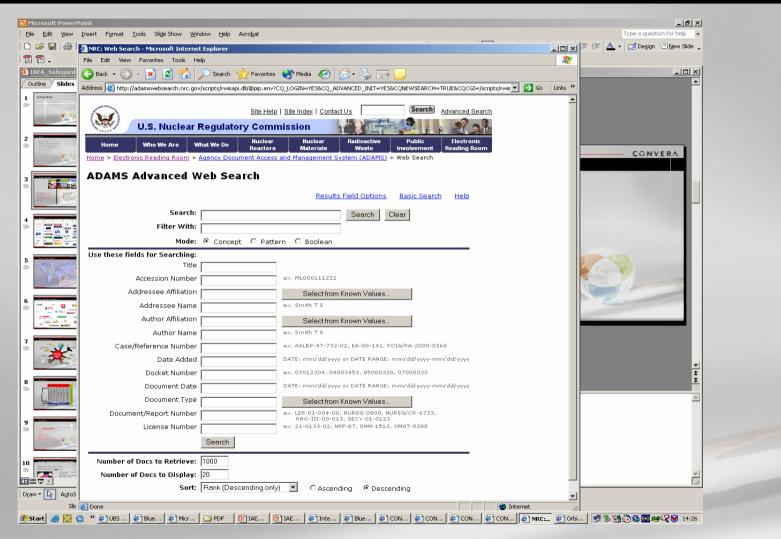# Convera & Partner Offices



**CONVERA SALES OFFICES**
**NORTH AMERICAN SALES PARTNERS**
**ASIA PACIFIC SALES PARTNERS**
**EUROPEAN SALES PARTNERS**
**ROW SALES PARTNERS**

*Convera Confidential*

CONVERA

## Our Focus

**LIFE SCIENCES**
- Johnson & Johnson
- Aventis
- NOVARTIS
- Pfizer
- AstraZeneca
- Roche
- gsk GlaxoSmithKline

**FINANCIAL**
- ABN·AMRO
- Goldman Sachs
- HSBC
- The Royal Bank of Scotland
- GREENWICH CAPITAL
- The Federal Reserve Board
- Bloomberg

**HIGH TECH**
- COGNOS
- Netegrity
- PTC
- intel
- intraspect
- UNISYS We have a head for e-business.
- symantec

**GOVERNMENT**
- Social Security Administration USA
- USGS science for a changing world
- NASA
- Sandia National Laboratories
- IRS
- FEDERAL TRADE COMMISSION
- FDA

- abc NEWS
- ESPN
- Discovery Channel, TLC, Discovery Health, Animal Planet, travel
- Chicago Tribune INTERNET EDITION
- NATIONAL GEOGRAPHIC TELEVISION
- THOMSON

*Convera Confidential*

**CONVERA**™

**Law enforcement**

**(Europe)**

Interpol

Polisen

www.polisen.se

**....and ?**

**Intelligence**

**?**

**Organizations**

**(Nuclear Industry)**

**UK AEA**

Sandia
National
Laboratories

*Convera Confidential*

# Agenda

- About CONVERA

- Introduction

- CONVERA Platform Infrastructure

- CONVERA's Retrieval Technologies

- Visualisation

- Automatic Categorization Technologies

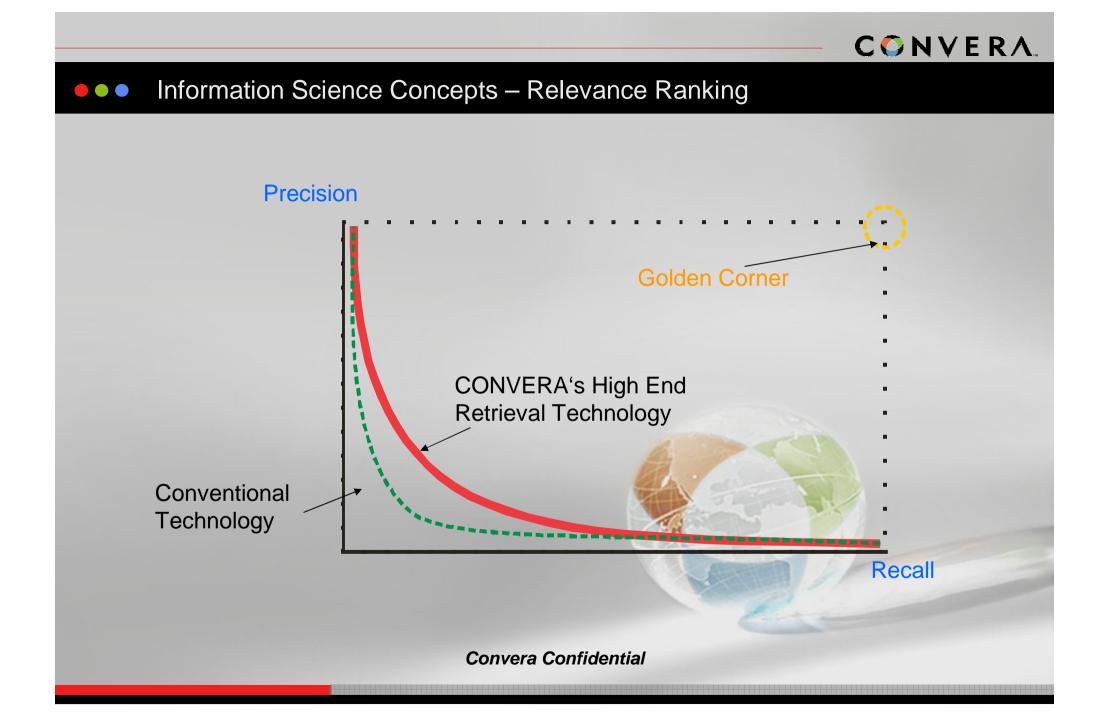- Dynamic Classification Technologies

- What we do / What's Third Party

- Product demonstration

*Convera Confidential*

Precision

Golden Corner

Recall

1. Completness
2. Proximity
3. Hit Density
4. Semantic Distance
5. Contextual Evidence

Precision

Golden Corner

CONVERA's High End
Retrieval Technology

Conventional
Technology

Recall

## Information Retrieval

1. We know what we know

   Example: Today's Weather

2. We know what we don't know

   Example: Winner of 1988 Oscars
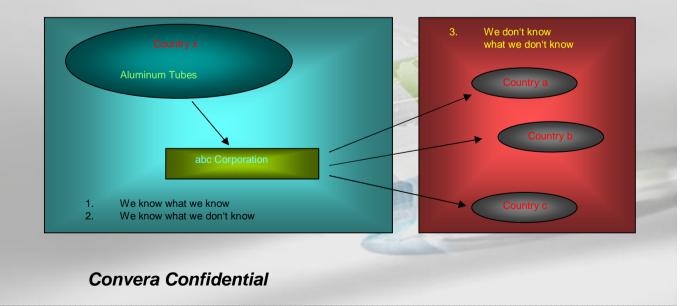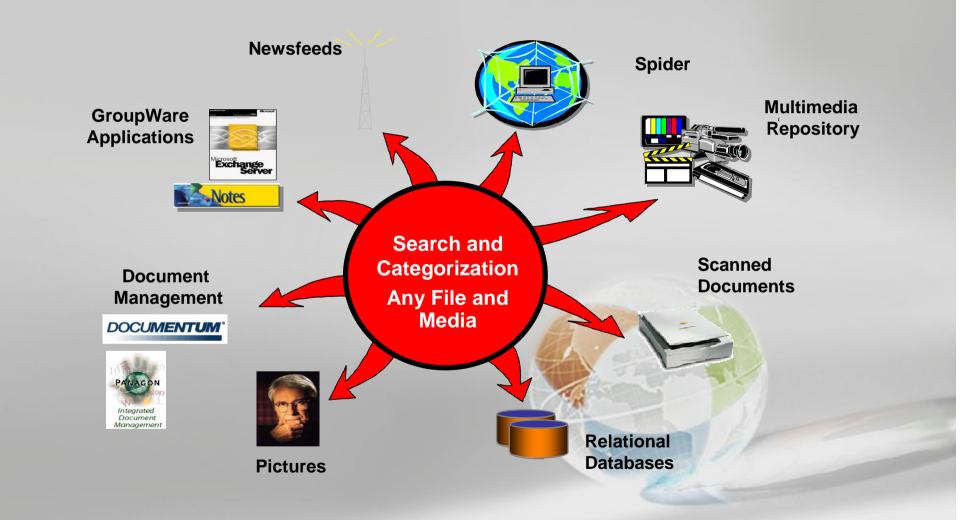
3. We don't know what we don't know

   Example: Discovering unknow user of dual purpose Technology

*Convera Confidential*

CONVERA

Country x

Aluminum Tubes

abc Corporation

1. We know what we know
2. We know what we don't know

3. We don't know
what we don't know

Country a

Country b

Country c

*Convera Confidential*

## Information Retrieval – How to do that?

1. Same as (Relevance Feedback)
2. Fuzzy Searching (Different Spelling, OCR Errors, ...)
3. Sematic Networks (Diff Corp. Names, Abreviations,
4. Taxonomies (Automatic Indexing)
5. Dynamic Classification (Directory is the Result – not the starting point)

Country x

Aluminum Tubes

abc Corporation

1. We know what we know
2. We know what we don't know

3. We don't know what we don't know

Country a

Country b

Country c

*Convera Confidential*

- About CONVERA

- Introduction

- CONVERA Platform Infrastructure

- CONVERA's Retrieval Technologies

- Visualisation

- Automatic Categorization Technologies

- Dynamic Classification Technologies

- What we do / What's Third Party

- Product demonstration

**Newsfeeds**

**Spider**

**GroupWare Applications**

**Multimedia Repository**

**Search and Categorization Any File and Media**

**Document Management**

**DOCUMENTUM**

**PANAGON**
*Integrated Document Management*

**Scanned Documents**

**Pictures**

**Relational Databases**

*Secure, standards-based integration for convenient, single-point of access.*

B6

# Agenda

- About CONVERA

- Introduction

- CONVERA Platform Infrastructure

- CONVERA's Retrieval Technologies

- Visualisation

- Automatic Categorization Technologies

- Dynamic Classification Technologies

- What we do / What's Third Party

- Product demonstration

*Convera Confidential*

## Technologies

- APRP – **Pattern recognition**
- List Search Server
- Boolean
- Concept
- Language Modules
- Multi Language Modules
- Cross  Language Modules
- Information Profiling
- Automatic Categorization
- Dynamic Personal Classification

*Convera Confidential*

1) Muammar Qaddafi
2) Mo'ammar Gadhafi
3) Muammar Kaddafi
4) Muammar Qadhafi
5) Moammar El Kadhafi
6) Muammar Gadafi
7) Mu'ammar al-Qadafi
8) Moamer El Kazzafi
9) Moamar al-Gaddafi
10) Mu'ammar Al Qathafi
11) Muammar Al Qathafi
12) Mo'ammar el-Gadhafi
13) Moamar El Kadhafi
14) Muammar al-Qadhafi
15) Mu'ammar al-Qadhdhafi
16) Mu'ammar Qadafi

17) Moamar Gaddafi
18) Mu'ammar Qadhdhafi
19) Muammar Khaddafi
20) Muammar al-Khaddafi
21) Mu'amar al-Kadafi
22) Muammar Ghaddafy
23) Muammar Ghadafi
24) Muammar Ghaddafi
25) Muamar Kaddafi
26) Muammar Quathafi
27) Mohammer Q'udafi
28) Muammar Gheddafi
29) Muamar Al-Kaddafi
30) Moammar Khadafy
31) Moammar Qudhafi
32) Mu'ammar al-Qaddafi

*Convera Confidential*

**BOOT** `01000010 01001111 01001111 01010100`

**BOAT** `01000010 01001111 01000001 01010100`



- **Overcomes errors, typos, misspellings**
- **25% inaccuracy in text is only 10% in binary**
- **APRP supports multimedia**

Sample: Names
Rayin al-Abidin Muamar Hussayn
Rayn al-Abidin Muhammar Husayn

**CONVERA**

- **Cross-lingual**
  - English
  - German
  - Spanish
  - French
  - Dutch
  - Italian
  - Hebrew and Arabic

*Allow stakeholders to access and view results across languages*

Bank

SEARCH

powered by **CONVERA**

**bank
S&L**

ENGLISH

**banque
coffre fort**

FRENCH

**depositar
caja fuerte**

SPANISH

*Convera Confidential*

# Applications of Concept Search

## Names and Aliases

Elizabeth

SEARCH powered by CONVERA

- **Elizabeth**
  - ➲ Lizzy — **ENGLISH**

- **Isabel**
  - ➲ Isabellita — **SPANISH**

## Organizations

Acme Holding Co

SEARCH powered by CONVERA

- **Acme Holding Company**
  - ➲ **Acme Widget, Inc.**
    - ➲ **Ohio Facility**
  - ➲ **Acme Import Export SA**
  - ➲ **Acme Shipping**

## Industry Terms

Stock

SEARCH powered by CONVERA

- **Stock**
  - ➲ **Securities**
  - ➲ **Equities**

*Provide stakeholders with 'virtual expertise' for more accurate search*

# Search Functionality – In the BODY TEXT and in the FIELDS!

- **Query by Example (relevance feedback)**
- **Idiom (Syntactic) Processing**
- **Adjustable Stop Words**
- **Exact Phrases**
- **Date Ranges**
- **Fielded Searching**
- **Search Term Weighting**
- **Logging functions (user)**
- **Web crawler**

- **Numeric Range searching**
- **Multiple Dictionaries / Thesauri**
- **Recurrent searching (searching hitlist)**
- **Multiple options for document display**
- **Automatic categorization**
- **Dynamic classification**
- **User profiling**
- **Relevance Ranking**
- **Language / Industry Plug-ins**
- **Cluster displays**
- **Visualisation aids**

Scalability  Accuracy  *Flexibility*

| Third Party | Visualization | Tracking | Pattern Detection |
|---|---|---|---|
| | Search | Classification | Mapping |
| CONVERA | Indexing | Categorization | Entity Extraction |
| | Crawling, Filtering | | |

- Clustering
- Statistical Analysis
- Absolute Linking
- Theme Grouping
- etc …

CONVERA



**Convera Confidential**

# Agenda

- About CONVERA

- Introduction

- CONVERA Platform Infrastructure

- CONVERA's Retrieval Technologies

- Visualisation

- Automatic Categorization Technologies

- Dynamic Classification Technologies

- What we do / What's Third Party

- Product demonstration

- Manual indexing is costly and slow

- Traditional Classification is precoordinated

- Hit lists are OK, but somehow inefficient

- Most information is unstructured

- Information structure is irrelevant

*Convera Confidential*

## Ontology

- An **ontology** is a foundation of categories representing a view of the world. An ontology reflects the commonly used and trusted breakdown of categories. For example, the breakdown of news items into categories of 'World', 'Sports', 'Politics', etc. is ontological.

## Taxonomy

- A **taxonomy** is a hierarchical system describing genera and species. Species derive from a common genus and are hierarchically represented according to their essential characteristics and differences. For example, animals are categorized with the "Taxonomy of Life" which separates mammals from birds and spiders from insects, based on proper features and relative differences. This genus to species nomenclature is highlighted by terminology which moves from generic terms to binomial terms through lexical derivation and compounding.

- A taxonomy doesn't deal with things, but with the essence of things: a taxonomy is based on an ontology.

# Categorization vs. Classification

- Categorization
  - **Logical**
  - **Taxonomy based**
  - **Consistent (Based on cultural fundamentals)**
  - **Stable**

- Classification
  - **Pragmatic**
  - **Precordinated**
  - **Common sense**
  - **Chaotic (Based on best practices)**
  - **<u>Individual</u>**

# Consistent, Scalable, and Flexible Knowledge Architecture

## We don't know the classification from tomorrows needs!

taxonomies

MeSH    Geography    INIS

classifications

Countries    Diseases    Communications

viewers

content

Integrated, Secure Indexing & Tagging

Index with Tags

tags

Fast, Relevant Classification

docs

- Category Browsing
- Dynamic Classification
- Visual Discovery

*Convera Confidential*

## Example 1: Geography

Africa
   Algeria
   Angola
Asia
   Afghanistan
   Armenia
Europe
   Albania
   Andorra
Middle East
   Bahrain
   Iran

North and Central America
   Antigua and Barbuda
   Bahamas
Pacific
   Australia
   Fiji
South America
   Argentina
   Bolivia
U.S.
   Alabama
   Alaska

# Example 2 : Defense

Defense Communications
  Satellite Communications
  Tactical Communications
Defense Systems
  Air Defense
          Antiaircraft Defense Systems
                  Gun Air Defense Systems
          Antimissile Defense Systems
          Forward Area Air Defense Systems
          Terminal Defense
  Aircraft Defense Systems
  Antisubmarine Defense Systems
  Antiswimmer Defense Systems
  Countermeasures
          Acoustic Countermeasures

Europe

Scandinavia

Finland

Schweden

Stockholm

Malmö

Plug and Play

Visual Discovery

Convera Confidential

## Expand Your Search

| | Sales | Marketing | Etc. |
|---|---|---|---|
| Compensation | | | |
| Career | | | |
| Etc. | | | |

- Search is a journey through a multi-dimensional grid of topics
- The ability to visualize all possible combinations at once will save time and increase focus
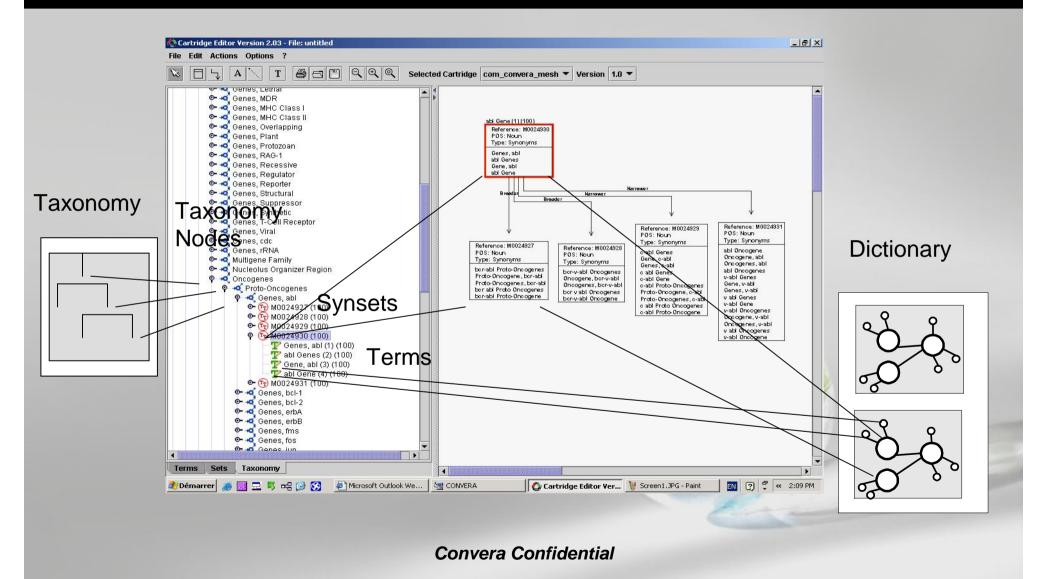
# Already Existing Taxonomy Cartridges (Samples)

- **Biology**
- **Chemistry**
- **Computers**
- **Electronics**
- **Finance**
- **Food Science**
- **Geography**
- **Geology**
- **Health Sciences**
- **Information Science**
- **Law**
- **Mathematics**

- **MeSH (Medical Subject Headings)**
- **Military**
- **Petroleum Natural Gas & Petrochemicals**
- **Pharmacology**
- **Physics**
- **Plastics**
- **Rubber**
- **Telecommunications**

Taxonomy

Taxonomy Nodes

Synsets

Terms

Dictionary

**Tables - Categorize**

| Category (Tax) / N Categories (Tags) | Document / N Categories | Category / N Latches | Category / N Ambiguou: |
|---|---|---|---|

|  | total | LOCAL- | depth | ndesc | category | path |
|---|---|---|---|---|---|---|
| 2 | 382 | 382 | 3 | 1 | Missiles | Ordnance, Missiles, and Munitions::Weapons::Missiles |
| 3 | 570 | 293 | 2 | 18 | Agents | Military Warfare::Agents |
| 4 | 333 | 287 | 3 | 3 | Terrorism | Military Warfare::Unconventional Warfare::Terrorism |
| 5 | 280 | 280 | 2 | 1 | National Security | Intelligence and Counterintelligence::National S |
| 6 | 630 | 269 | 1 | 16 | Intelligence and Counterintelligence | Intelligence and Counterint |
| 7 | 231 | 228 | 2 | 5 | Missions | Military Operations::Missions |
| 8 | 229 | 211 | 2 | 30 | Bombs | Ordnance, Missiles, and Munitions::Bombs |
| 9 | 209 | 209 | 2 | 4 | Battles | Military Warfare::Battles |
| 10 | 260 | 198 | 3 | 5 | Anthrax | Military Warfare::Agents::Anthrax |
| 11 | 815 | 174 | 2 | 58 | Weapons | Ordnance, Missiles, and Munitions::Weapons |
| 12 | 169 | 169 | 3 | 1 | Guns | Ordnance, Missiles, and Munitions::Chemical Ordnance::C |
| 13 | 162 | 160 | 1 | 3 | Defense Communications | Defense Communications |
| 14 | 96 | 75 | 3 | 7 | Nuclear Weapons | Ordnance, Missiles, and Munitions::Weapons::N |
| 15 | 58 | 58 | 3 | 1 | Peacekeeping | Military Operations::Operations Other Than Wa |
| 16 | 396 | 46 | 1 | 34 | Military Operations | Military Operations |
| 17 | 46 | 46 | 3 | 1 | Biological Weapons | Ordnance, Missiles, and Munitions::Weapons::E |
| 18 | 36 | 36 | 3 | 1 | Persian Gulf War | Military Warfare::Chemical Warfare::Persian Gu |
| 19 | 36 | 36 | 4 | 2 | Rifles | Ordnance, Missiles, and Munitions::Chemical Ordnance::S |
| 20 | 35 | 35 | 3 | 2 | Grenades | Ordnance, Missiles, and Munitions::Ammunition::Grenade: |
| 21 | 35 | 35 | 4 | 1 | Assassination | Military Warfare::Unconventional Warfare::Terr |

# Charts to improve the Cartridge Quality

**CONVERA**

Relate available information to

YOUR

decision-making processes

- Categorize with consistency
- Classify in context

*Convera Confidential*

- **The directory is a result, not a starting point.**

- **Ontologies are real ontologies: conceptual and explicit.**

## Typical Structures

- Geography / Topic
  - **Terrorism in Philippines**
  - **Criminal Law in Texas**
  - **Domestic Sales**
  - **Security in Building C**

- Horizontal / Vertical
  - **Petroleum Business**
  - **AML Regulations**

- Vertical / Vertical
  - **Chemical Compounds for Alzheimer**

## Over Defined Context

- Very large computational space
  - **"Chemical Compounds in Alzheimer Genomics"**
  - **-> 8500 diseases**
  - **-> 1,000 genes**
  - **-> 30,000 compounds**
  - **= 255 billion folders**

- Reversely proportional number of successfully populated folders

- Can't be done by automatic CLASSIFICATION!

## What We Do (Some Samples)

- List Search (Batch Mode)

- Information Profiling

- Automatic Categorizing

- Dynamic Classification

- Multi- and Cross Language

- Content Management

- Voice To Text and Automatic Meta Data Generation

## Third Party (Some Samples)

- Pre Processing (nCase, etc.)

- Post Processing (Statistics, Facerec, …)

# RetrievalWare 8
## High End Categorization/Classification

Heinz Bachmann