DEFENCE **R&D** DÉFENSE

# CAPTURING AND MODELING DOMAIN KNOWLEDGE USING NATURAL LANGUAGE PROCESSING TECHNIQUES

Alain Auger, Ph. D.

IKM Section / DRDC Valcartier

June 2005

(Paper 296)

Defence Research and Development Canada

Recherche et développement pour la défense Canada

Canadä

# Problem Space

- Command and control (C2) and decision-making domains are seriously threatened facing information overload and uncertainty issues

- Military have to create new ways of processing sensor and intelligence information

- **Without new means to elicit knowledge from multiple information and intelligence sources, decision-makers will have to deal with very limited knowledge and increasing levels of uncertainty in operations**

- How can we better capture and represent knowledge objects contained in sources?

# Knowledge Representation Enablers

- Metadata

- Taxonomies

- Ontologies

# Some Metadata Sets

- **Metadata** (Greek: *meta-* + *data* "information") means « data about data ».

- Dublin Core

  – The **Dublin Core Metadata Element Set** consists of 16 optional metadata elements, any of which may be repeated or omitted. (Title, Creator, Subject, Description, Publisher, Contributor, Date, Type, Format, Identifier, Source, Language, Relation, Coverage, Rights, Audience)

- Resource Description Framework (RDF)

  – The purpose of RDF is to provide an encoding and interpretation mechanism so that resources can be described in a way that particular software can understand it, or, better put, so that software can more easily access data organized within structured parameters.

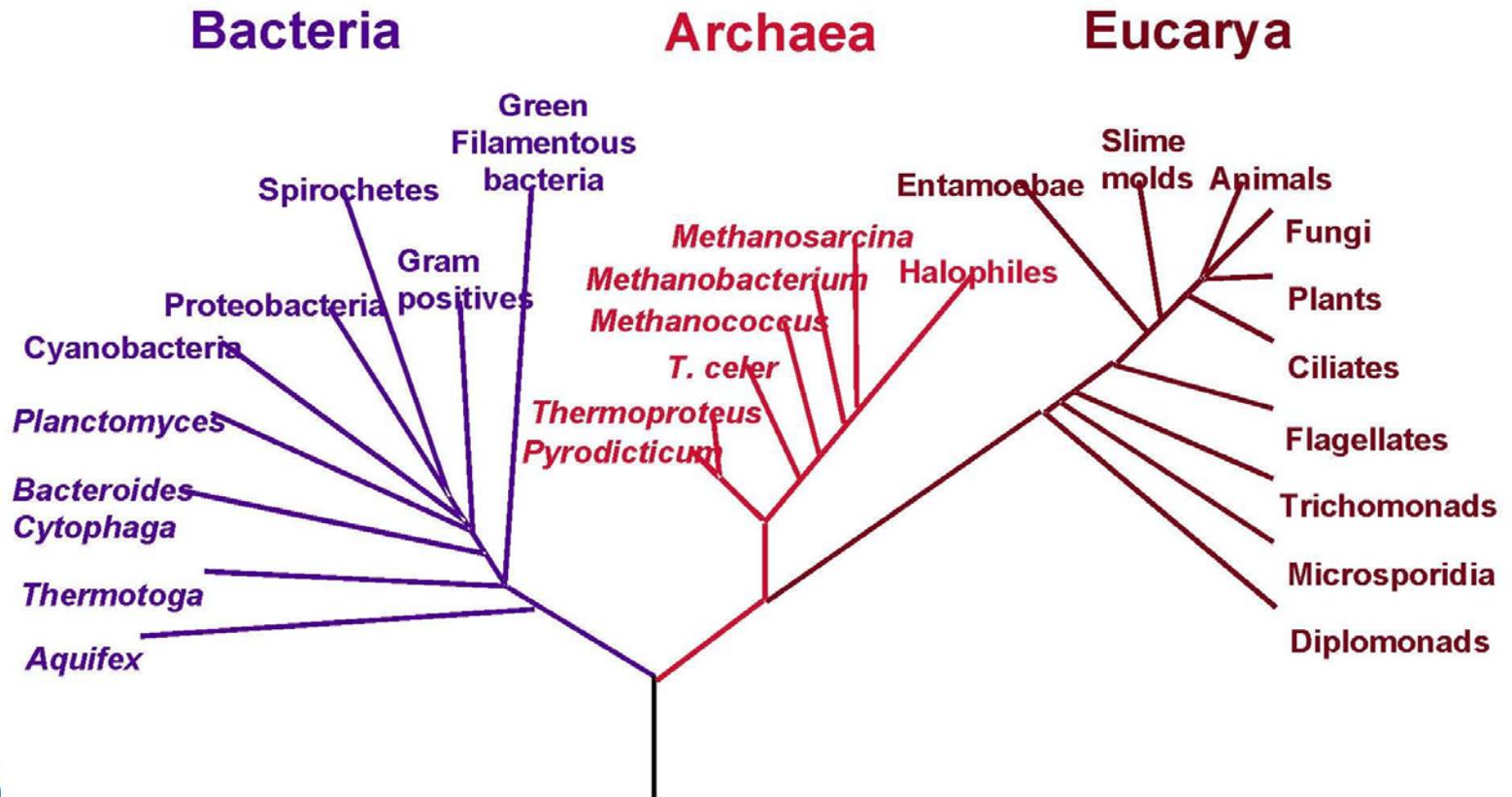- Extended Markup Language (XML)

- Etc.

# Taxonomies

- Taxonomy (from Greek ταξινομία (*taxinomia*) from the words *taxis* "order" and *nomos* "law") may refer to either a hierarchical classification of things, or the principles underlying the classification. Almost anything, animate objects, inanimate objects, places, and events, may be classified according to some taxonomic scheme. [Wikipedia]

- In taxonomies, concepts are classified using **homology**; that is, **shared characteristics that have been inherited from a common ancestor**.

- **Limitation**: IS-A or PARENT-CHILD relationship type only. Cannot express CAUSE-EFFECT relationships, for instance
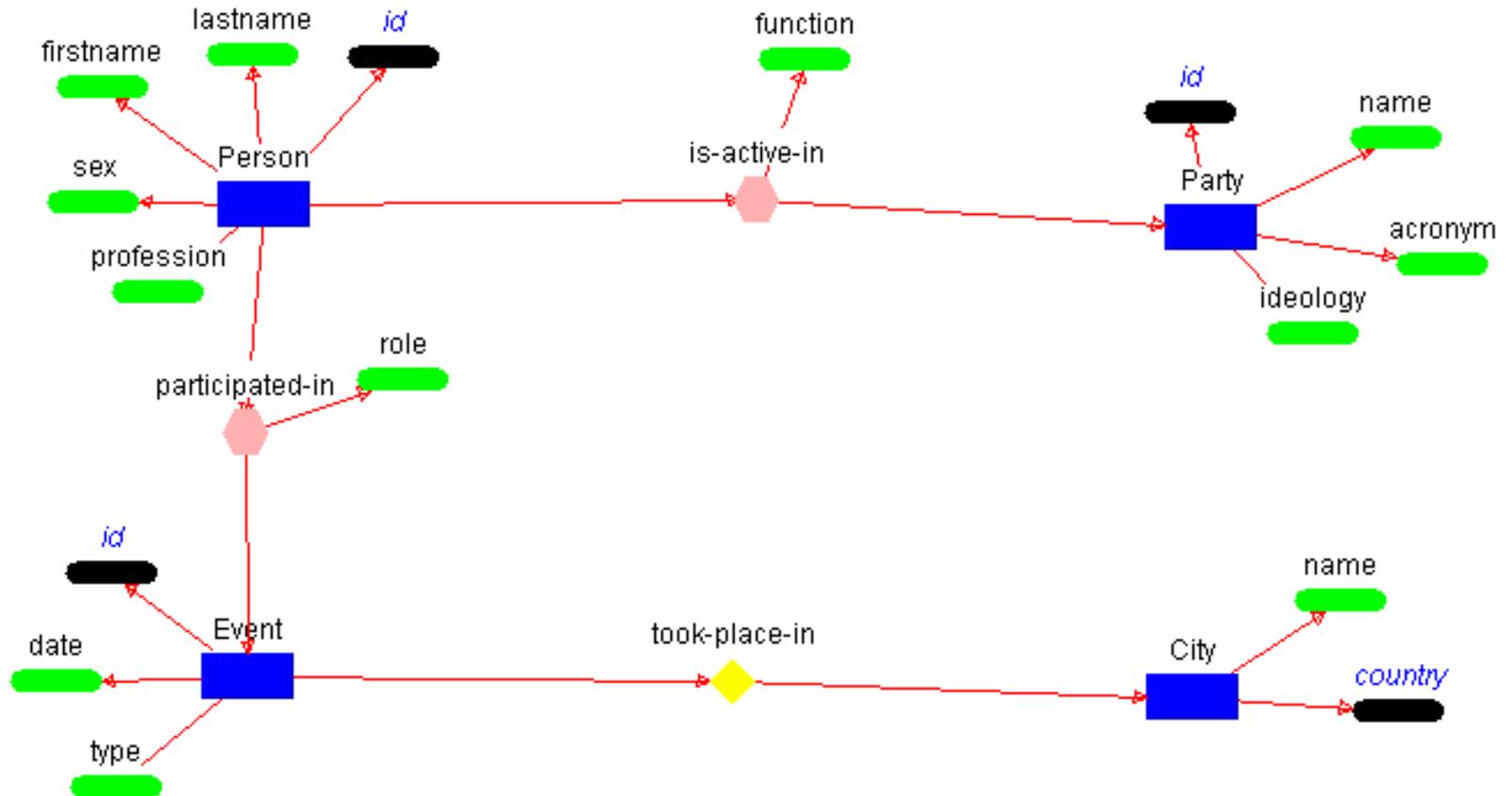
# Taxonomy Sample

# Phylogenetic Tree of Life

# Ontologies

- An **ontology** is a formal, explicit specification of a shared conceptualisation [Gruber, 1993]

- An ontology is a formal explicit specification of how to represent the objects, concepts, and other entities that are assumed to exist in some area of interest and the relationships that hold among them.
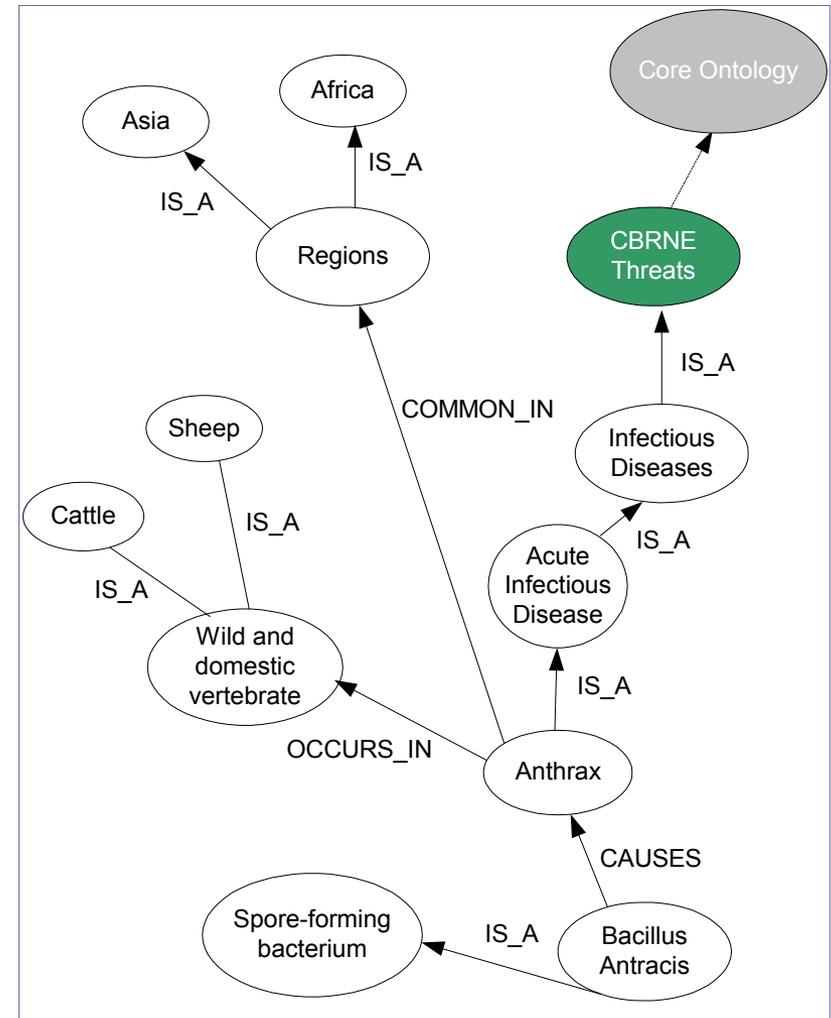
# Ontology Sample

# The Need for Domain Ontologies

- Domain ontologies are key elements required to enable next generation of decision support and knowledge exploitation systems with new semantic capabilities

- Ontology-engineering remains a non-trivial, time and budget consuming activity

- How can we rapidly build ontologies?

# SACOT Research Project

- Aim:

  – To develop and apply natural language processing (NLP) extraction techniques to unstructured texts to capture knowledge objects they contain and represent them in the form of an ontology

# Limitations of Traditional Ontology-Engineering Approaches

- Relying on Humans

    - Based on Subject Matter Experts

    - Adapted to task or application ontologies

    - Not adapted to domains ontologies (too many objects)

- Relying on Statistics

    - e.g. computation of co-occurring words

# SACOT Ontology-engineering Process

- Sources Identification

- Extraction Processes

- Draft Ontologies Generation

- Draft Ontologies Validation

- Ontology Maintenance

# What are Domain Ontologies Made of?

- Named Entities expressed in texts

- Concepts expressed by terms in texts

- Relations among knowledge objects

# SACOT's Specifics

- Domain-specific Named Entity Extraction

- Contrastive Approach to Terminology Extraction

- Natural Language Processing (NLP) approach to semantic relations extraction

# Named Entities Extraction

# Terminology Extraction

- SACOT's Specifics:
  - Use of a contrastive approach to compute and automate candidate terms validation process

| Frequence | Term | Score |
|---|---|---|
| 6619 | terrorist | 101,99 |
| 4209 | terrorism | 92,80 |
| 4587 | nuclear | 83,01 |
| 3018 | biological | 78,67 |
| 2520 | weapon | 68,01 |
| 1895 | Iraq | 61,35 |
| 2107 | attack | 57,79 |
| 1885 | domestic | 55,80 |
| 1200 | department | 47,57 |
| 1125 | al | 47,18 |
| 2266 | military | 46,97 |
| 1527 | September | 46,59 |
| 1048 | Iraqi | 46,23 |

# Semantic Relations Extraction

# SACOT's Specifics

- Targeting specific semantic relations markers that are present in texts as explicit « indicators » to capture relations among concepts

  – e.g. *X is used to Y*, *X is located in Y*

- Not based on co-occurrence statistics

- Entirely based on semantic relation patterns

  – e.g. *is used to*, *is located in*

# Putting it All Together

**Sample Input Text**

May 24, 2002
Anthrax is an acute infectious disease caused by the spore-forming bacterium Bacillus anthracis. Anthrax most commonly occurs in wild and domestic lower vertebrates (cattle, sheep, goats, camels, antelopes, and other herbivores), but it can also occur in humans when they are exposed to infected animals or tissue from infected animals.

Anthrax is most common in agricultural regions where it occurs in animals. These include South and Central America, Southern and Eastern Europe, Asia, Africa, the Caribbean, and the Middle East. When anthrax affects humans, it is usually due to an occupational exposure to infected animals or their products. Workers who are exposed to dead animals and animal products from other countries where anthrax is more common may become infected with B. anthracis (industrial anthrax). Anthrax in wild livestock has occurred in the United States.

Automatic Terminology Extraction Process

Automatic Named Entities Extraction Process

Automatic Semantic Relations Extraction Process

**Candidate Terms**

anthrax
acute infectious disease
spore-forming bacterium
Bacillus anthracis
wild and domestic lower
vertebrates
cattle
sheep
goat

**Candidate Named Entities**

DATE: May 24 2002
GEONAME: South and
Central America
GEONAME: Southern and
Eastern Europe
GEONAME: Asia
GEONAME: Africa
GEONAME: Caribbean
GEONAME: Middle East
GEONAME: United States

**Candidate Semantic Relations**

anthrax IS_A acute infectious disease

Bacillus anthracis CAUSES anthrax

anthrax OCCURS_IN wild and domestic lower vertebrate

cattle IS_A wild and lower vertebrate

sheep IS_A wild and lower vertebrate

Validation

Validation
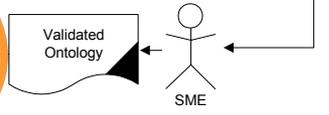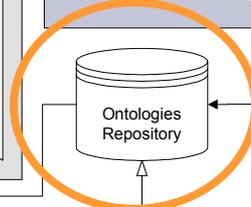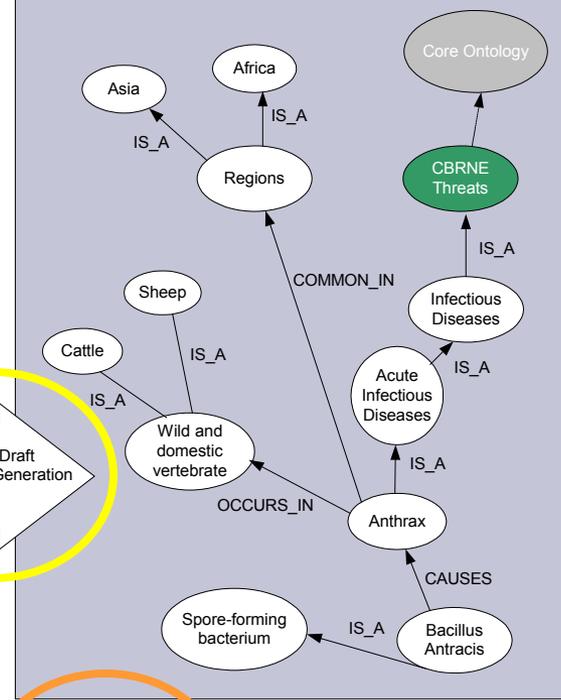
Validation

**Validated Lists**

anthrax
acute infectious disease
spore-forming bacterium
Bacillus anthracis
wild and domestic lower
vertebrates
cattle
sheep
goat

DATE: May 24 2002
GEONAME: South
America
GEONAME:  Central
America
GEONAME: Southern
Europe
GEONAME: Eastern
Europe
GEONAME: Asia
GEONAME: Africa
GEONAME: Caribbean
[...]

anthrax IS_A acute infectious disease

Bacillus anthracis CAUSES anthrax

anthrax OCCURS_IN wild and domestic lower vertebrate

cattle IS_A wild and lower vertebrate

sheep IS_A wild and lower vertebrate

Automatic Draft Ontology Generation Process

**Ontology Hypothesis**
*(to be validated by the SME)*

Core Ontology

Africa

Asia

Regions

CBRNE Threats

IS_A

IS_A

IS_A

COMMON_IN

IS_A

Infectious Diseases

Sheep

Cattle

IS_A

IS_A

IS_A

Acute Infectious Diseases

Wild and domestic vertebrate

OCCURS_IN

IS_A

Anthrax

CAUSES

Spore-forming bacterium

IS_A

Bacillus Antracis

Ontologies Repository

Validated Ontology

SME

**Ontology Services**

*Thrid Party Application (e.g. Knowledge Portal)*

# Conclusion

- Preliminary results show that the SACOT ontology-engineering framework might significantly reduces time usually required to capture the knowledge objects of a domain in traditional, fully human-based, ontology building processes.

# Project Status

- Initiated in 2004, SACOT is a research project in its early stage.

- All extraction modules are still under development

- All existing modules are standalone at the moment. They are not integrated in the SACOT framework.

# Way Ahead

- Measure performance of all three extraction modules

- Integrate all extraction modules in the SACOT framework

- Investigate machine learning techniques in support to SME validation of draft ontologies generated by the SACOT framework

DEFENCE **R&D** DÉFENSE