*Co-evolving document collections and knowledge structures*

# CoDAK

Dr. Evgeny Knutov

(MSc Seminar Nov. 11 2013)

# The CoDAK project

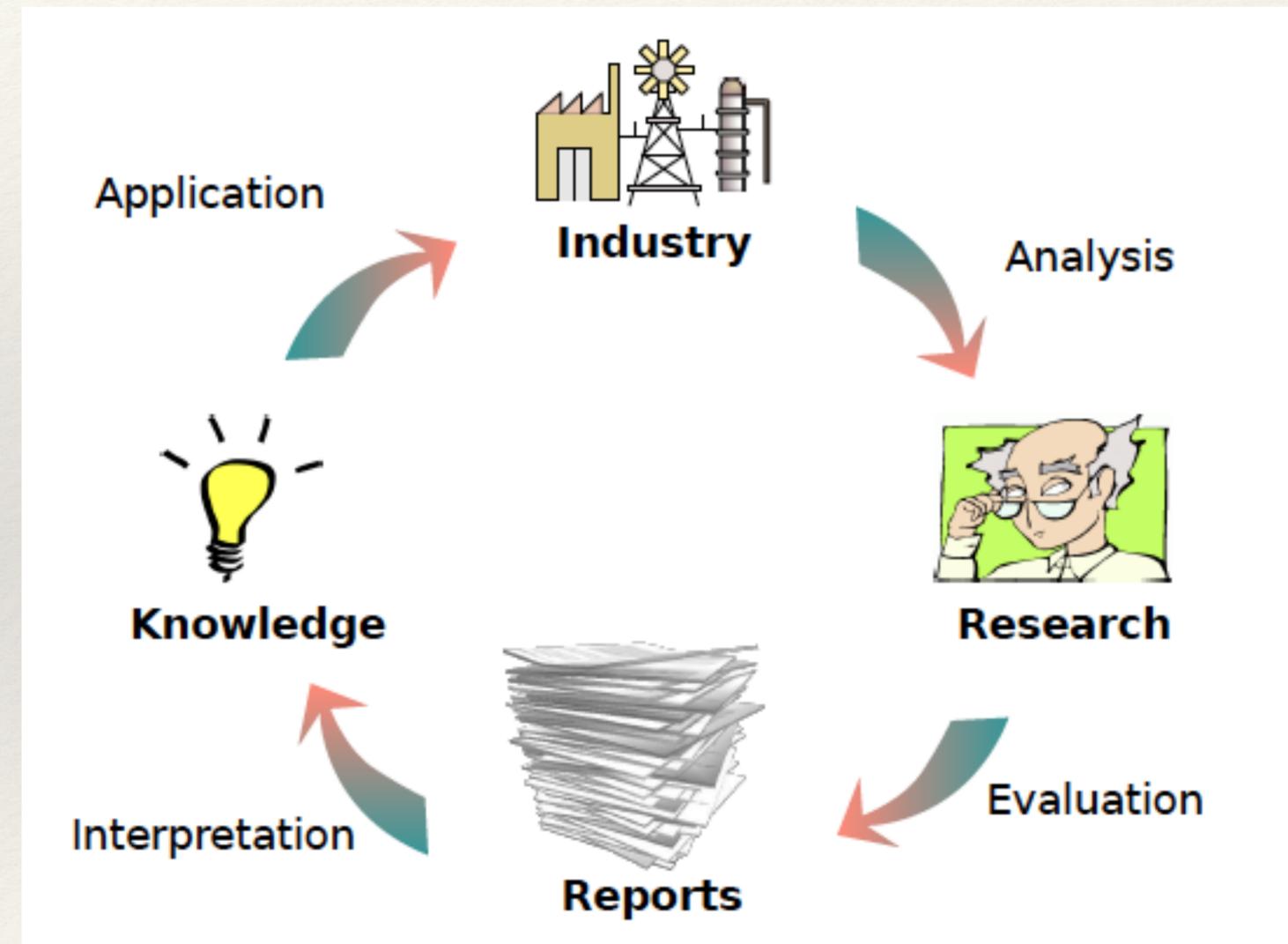CoDAK: Co-evolving Document Collections and Knowledge Structures

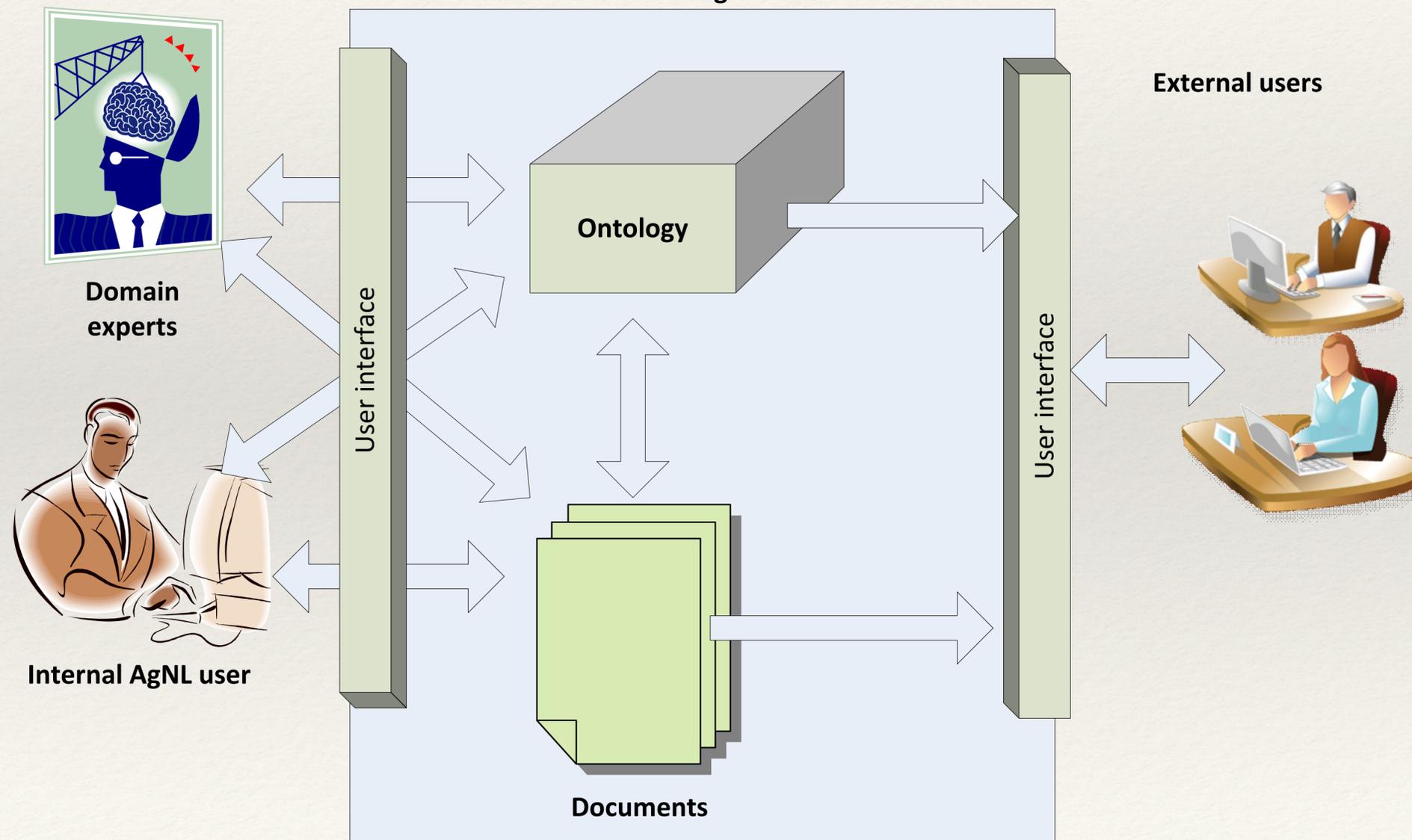AgentschapNL: dept. Energy & Climate http://www.agentschapnl.nl/

www.semontoweb.com

# The CoDAK project (cont.)

* Problem statement

* Research

* Evolution of CoDAK

* Demos

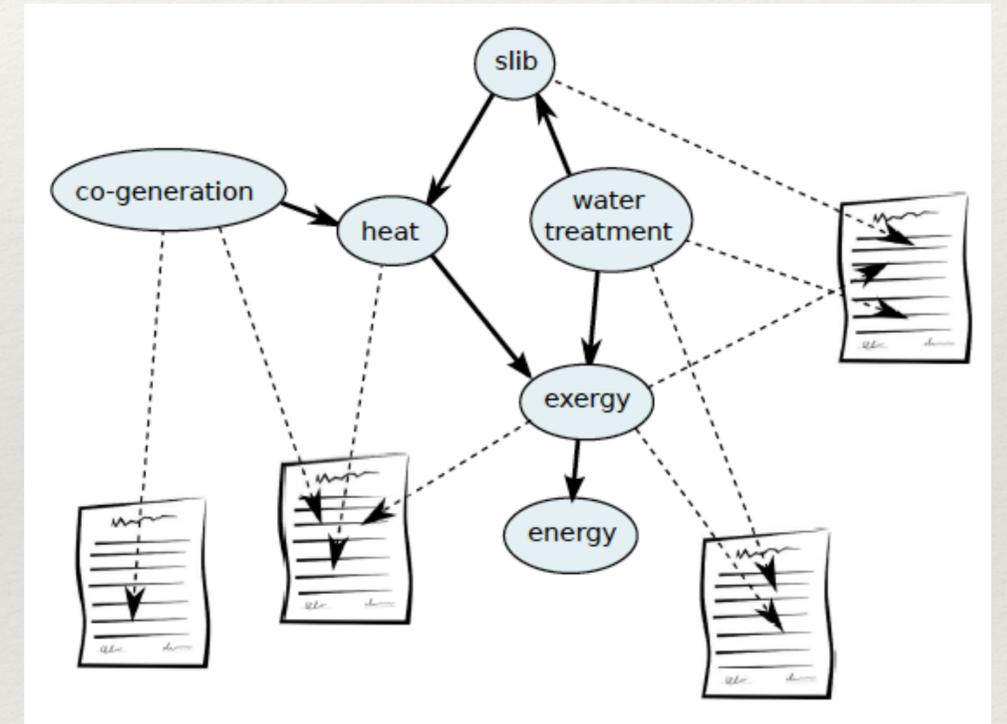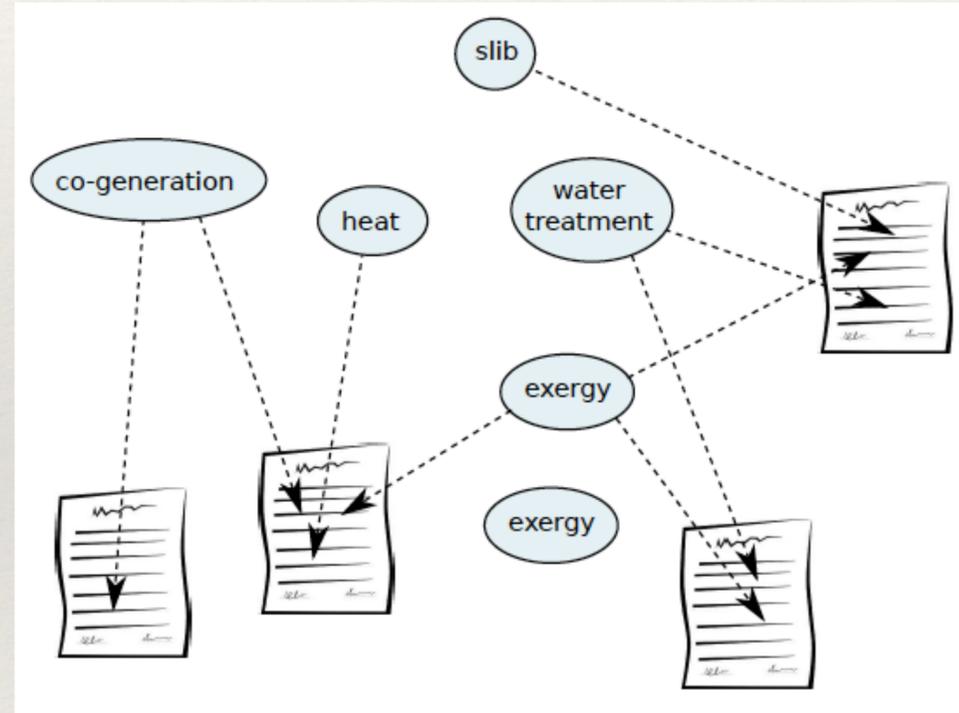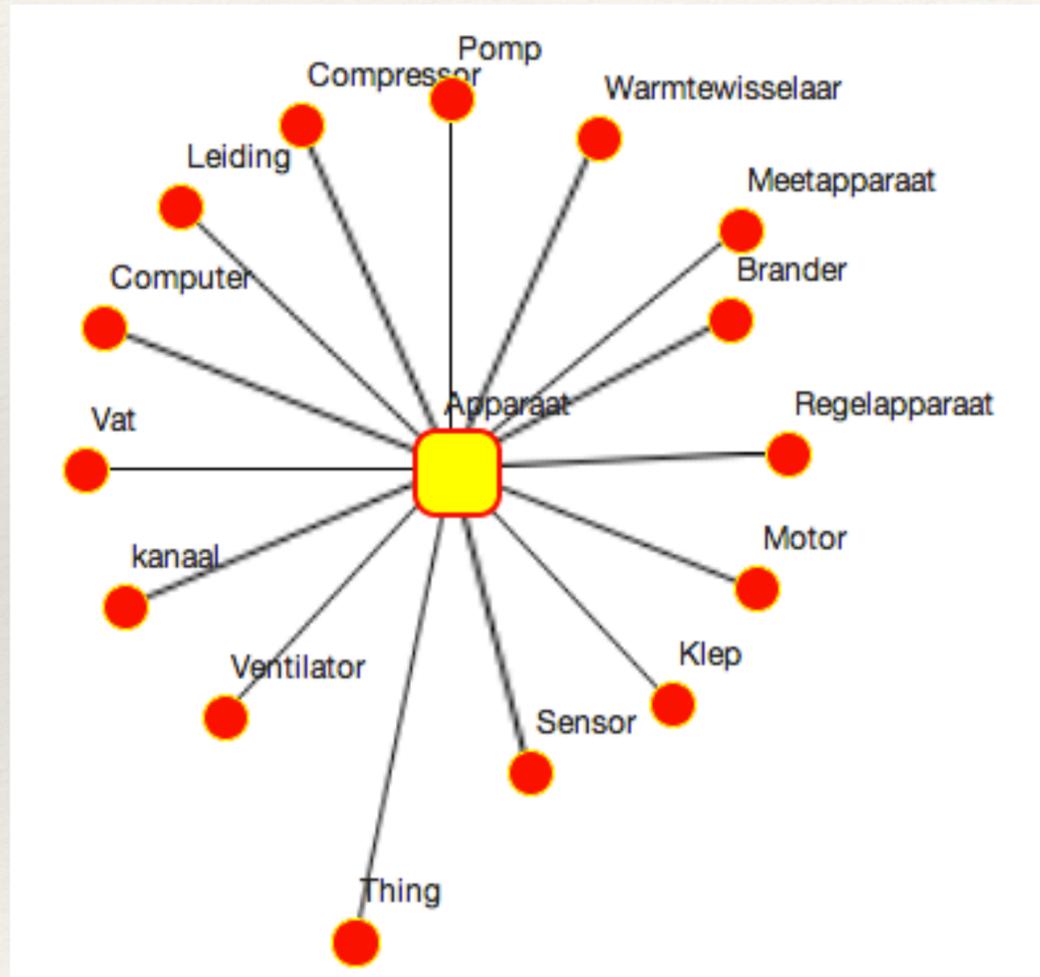* CoDAK language

* New research objectives

# How people at AgNL see the system



Dynamical system, knowledge evolves through new documents, creation of factsheets and further domain expert knowledge.

Domain experts

Internal AgNL user

User interface

Ontology

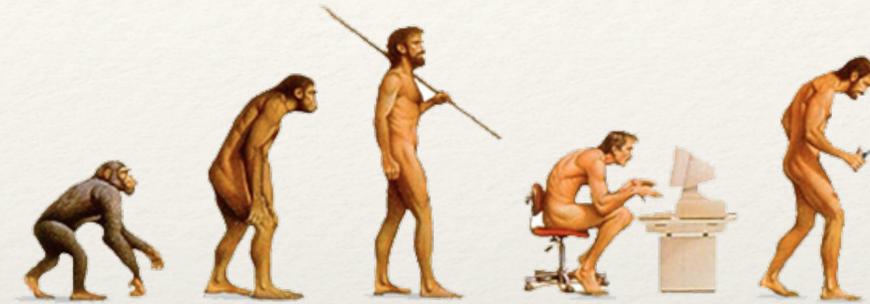Documents

User interface

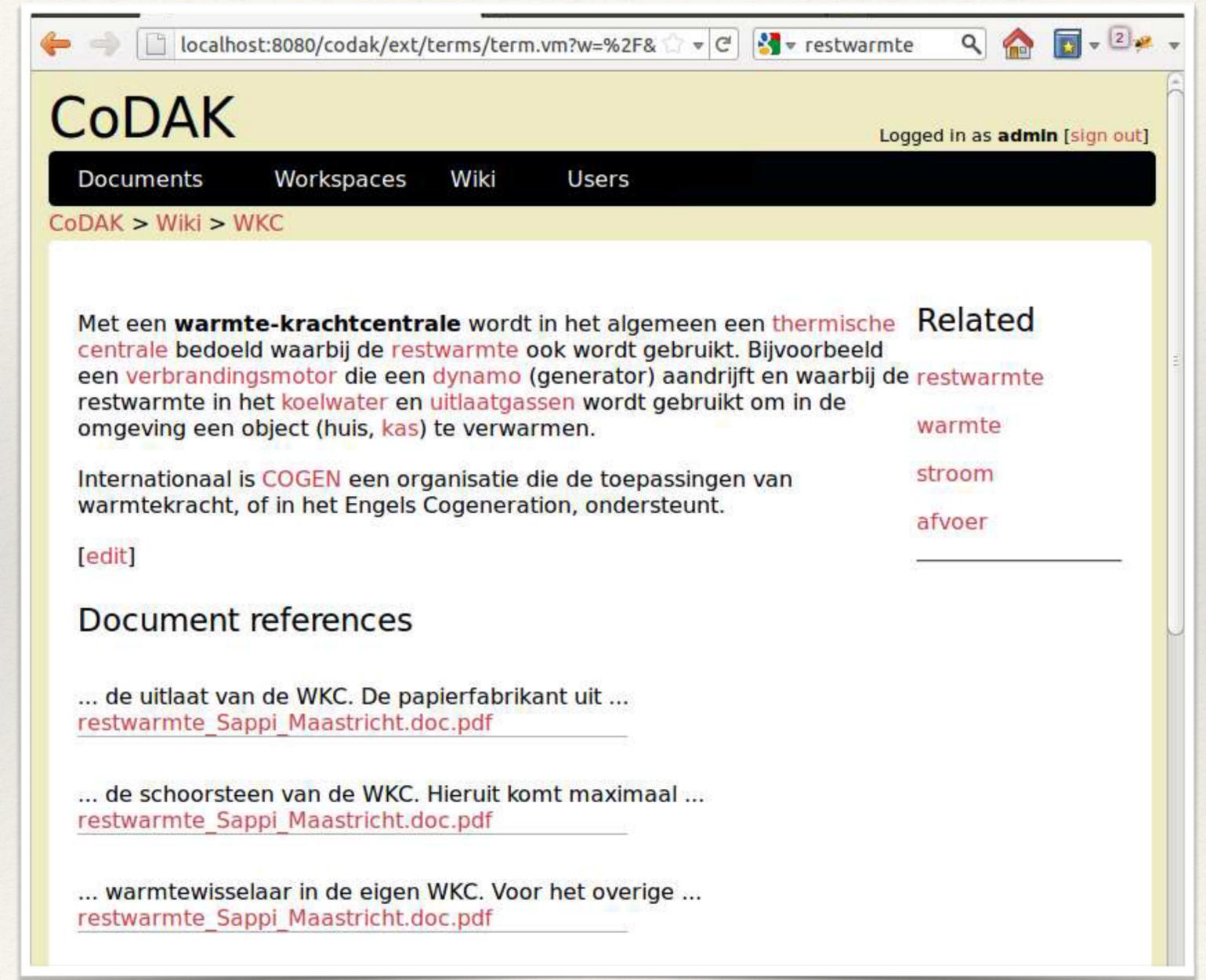External users

# Ontology and Documents

# Evolution of CoDAK

- ❖ multiple versions
- ❖ scrum approach
- ❖ deployment
- ❖ user studies
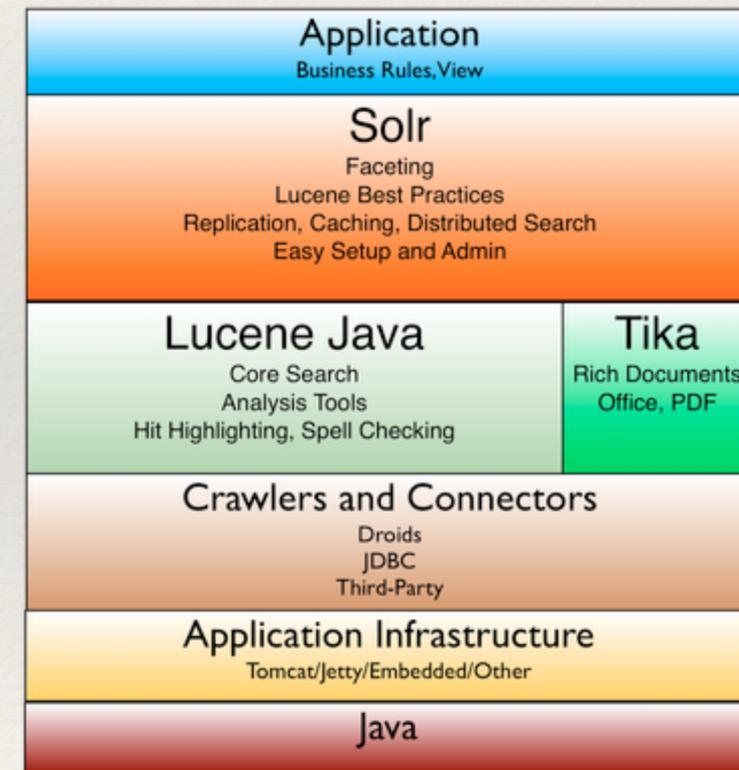- ❖ new endeavours

# First versions

- ❖ custom development

- ❖ uses Frog NLP module

- ❖ PDFBox for PDF documents processing

# Java tryouts

The Lucene Stack

- pure Java-based implementation

- includes:
  - **Lucene** indexed PDF files (with PDFbox)
  - **Sail** ontology index

- using Lucene Multi-index search



Application
Business Rules, View

Solr
Faceting
Lucene Best Practices
Replication, Caching, Distributed Search
Easy Setup and Admin

Lucene Java
Core Search
Analysis Tools
Hit Highlighting, Spell Checking

Tika
Rich Documents
Office, PDF

Crawlers and Connectors
Droids
JDBC
Third-Party

Application Infrastructure
Tomcat/Jetty/Embedded/Other

Java

# Drupal-based

- need to manage and create documents (friendly)

- Drupal CMS

- working towards "**WorkingSheets**"

- uses both internal search index (only for specific nodes) and external **Solr**

# Drupal ?

- open-source content management system (framework)

- consists of the core and additional modules (primarily PHP)

- data construction (types, views, versions)

- https://drupal.org/

- not?-out-of-the-box

- want something else to work together - do it yourself

- bunch of CGI

- new research challenges

  - FactSheet

  - UI

# Extracting for Drupal



extracted relationships

extracted entities

- ❖ Alchemy entity extraction
- ❖ Drupal taxonomies
- ❖ Automatic tagging

- ❖ nothing works with Drupal as expected;
- ❖ most of the Drupal modules are tailored to fit in specific implementation;
- ❖ lots of Ajax was failing in IE
- ❖ external API

# Analyzing content in Drupal



- automatically get the most relevant information of your report, table, etc. connected to many other documents and people to help you explore the domain;
  - identifies people, companies, geo features
  - automatically tags analyzed content
  - identifies facts and subject-object-action relationships
  - topic categorization

## Sustainable Water Fund: First call results

**Number of applications for the Sustainable Water Fund exceeds all expectations.**

On May 7th the first stage for the tender for the Sustainable Water Fund closed. A total of 81 applications was received. This exceeded even the most optimistic estimates. To the Ministry of Foreign Affairs and NL Agency this shows that there is great momentum for PPPs in the water sector and that the Sustainable Water Fund has the potential to really contribute to water safety and water security in developing countries.

The total requested indicative subsidy from the Fund is EUR 200 million, more than 4 times the available budget for this tender. This means that the Fund has managed to leverage roughly EUR 150 million from water sector actors in the Netherlands and abroad. The financial limitations of the tender, however will determine how many proposals can be funded. The Ministry of Foreign Affairs and NL Agency will keep the sector informed about this process.

Applications have been received for 27 target countries.

Read here more about the first call of result.

DELEN

Published on: 30-05-2012 | Changed on: 30-05-2012

Printer-friendly version

More about: Watersector, Afghanistan, Albanië, Armenië, Bangladesh, Benin, Bolivia, Bosnië en Herzegovina, Burkina Faso, Burundi, Colombia, Congo-Kinshasa, Egypte, Ethiopië, Filipijnen, Georgië, Ghana, Guatemala, Indonesië, Jemen, Kaapverdië, Kenia, Kosovo, Macedonië, Malawi, Mali, Marokko, Moldavië, Mongolië, Mozambique, Nicaragua, Oeganda, Pakistan, Palestijnse gebieden, Peru, Rwanda, Senegal, Sri Lanka, Suriname, Tanzania, Thailand, Vietnam, Zambia, Zuid-Afrika, Zuid-Soedan, Fonds Duurzaam Water

**More like this**

> Orio: New country list
> New ORIO Policy rules officially published
> The Facility for Infrastructure Development (ORIO)
> ORIO Call for Proposals 2012 officially opened

+ show more
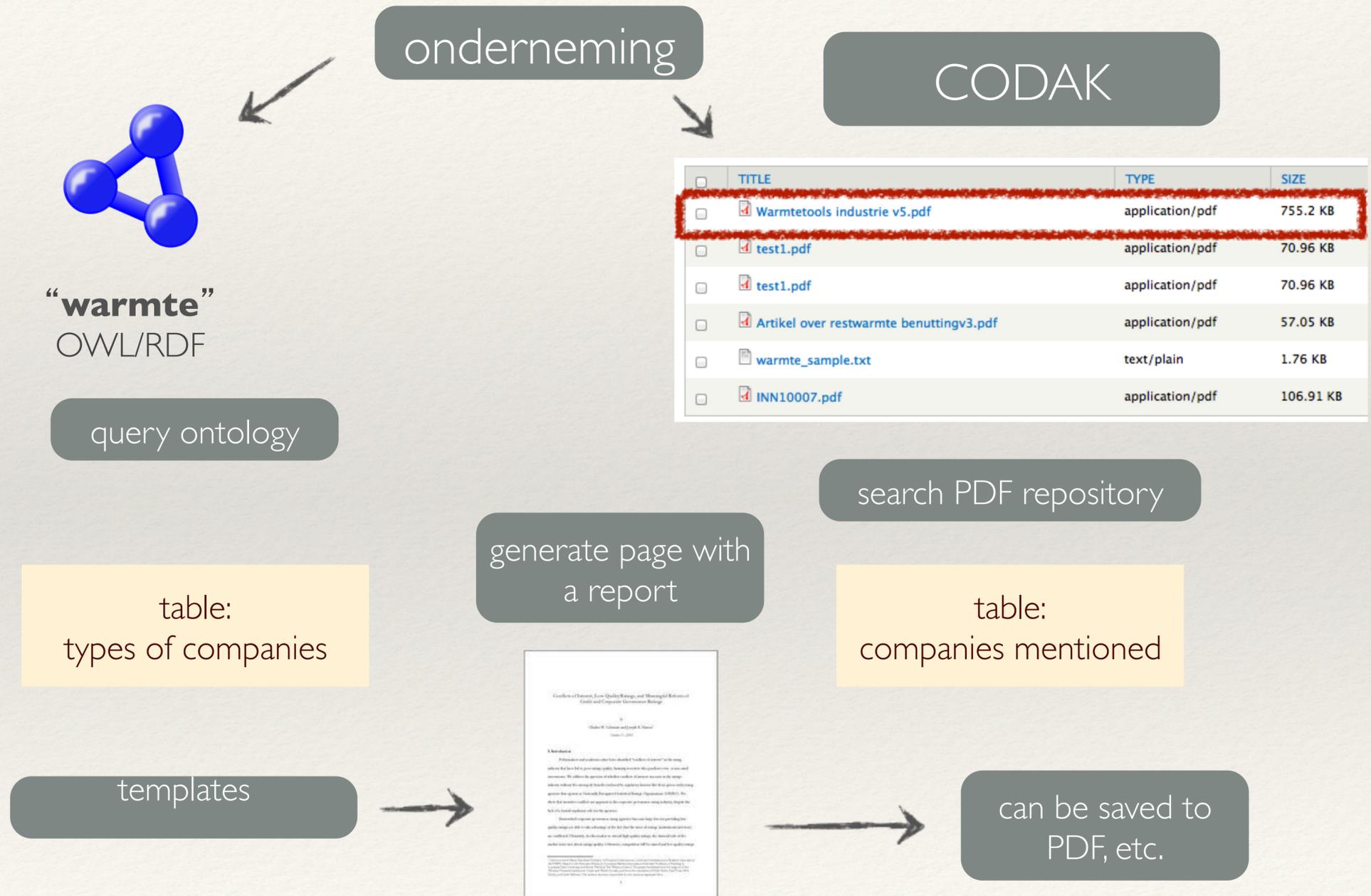
Agentschap NL
Ministerie van Economische Zaken,
Landbouw en Innovatie

**Actual
AgentschapNL web site**

# "Factsheet" generation tryouts

onderneming

CODAK

- ❖ generate simple reports
  - ❖ on a certain pre-defined topic
  - ❖ term-based

"**warmte**"
OWL/RDF

- ❖ using different report templates

- ❖ add reports to the document base

query ontology

| | TITLE | TYPE | SIZE |
|---|---|---|---|
| ☐ | Warmtetools industrie v5.pdf | application/pdf | 755.2 KB |
| ☐ | test1.pdf | application/pdf | 70.96 KB |
| ☐ | test1.pdf | application/pdf | 70.96 KB |
| ☐ | Artikel over restwarmte benuttingv3.pdf | application/pdf | 57.05 KB |
| ☐ | warmte_sample.txt | text/plain | 1.76 KB |
| ☐ | INN10007.pdf | application/pdf | 106.91 KB |

search PDF repository

generate page with a report

table:
types of companies

table:
companies mentioned

templates

can be saved to PDF, etc.

# "Factsheet" generation tryouts (cont.)



**CoDAK HUB**

Explore Warmte Ontology (jOWL) Query Ontology Explore Ontology (QPLE) Test CGI Run Ontology Visualizer Show Warmte hierarchy Show Warmte Graph

Generate Report by keyword [onderneming] [Submit]

**Report term**

## Onderneming

Een bedrijf is een organisatie van arbeid en kapitaal. Een bedrijf dat gericht is op het maken van winst wordt veelal een onderneming genoemd. Een bedrijf dat tastbare product
genoemd.
Bij commerciÃ«le organisaties is er sprake van een onderneming waar risico's worden gelopen en waar het streven naar winst noodzakelijk is voor het voortbestaan. Om hun kla
bereiken en aan te spreken met hun diensten, wordt er nauwkeurig gelet op de vier -P's-: prijs, plaats, product en promotie. Tegenwoordig handelen bedrijven ook steeds vaker v
natuurlijk Profit. In deze ruimere benadering wordt er ook rekening gehouden met de meerwaarde of kosten voor mens, maatschappij en de omgeving (m.a.w. duurzaam onder

**Template**

wiki information on this topic: onderneming  **wiki info**

And this is what we found in the 'warmte' Ontology related to :
Diendverlener - http://www.semontoweb.org/2011/warmte.owl#dienstverlener
Fabrikant - http://www.semontoweb.org/2011/warmte.owl#fabrikant

**Ontological information**

The following PDF files were found related to the topic of the report:
Vertrouwelijk eindrapport_IWBH10035.pdf

**Where in available PDFs**

# HTML5 + Graphics



❖ http://localhost/codak01/codak_10o.html

# HTML5 + Graphics (cont)

❖ badly scalable search engine (good UI though)

❖ lots of documents pre-processing (into HTML)

❖ buggy jOWL package

❖ IE5,6,7 visualization limitations

❖ limited set-up

# And again…, robust, scalable, simple HTML

❖ from scratch

❖ only search, no UI at first

❖ incrementally adding features

❖ textual representation for the ontology

❖ tracking users' behaviour

# First deployed versions based on

- **Solr,** starting from ver. 3.6

- current **Solr** ver. 4.5.1

- **Sesame** - Data Store and SPARQL backend

- using only XMLHttpRequest

- XML data only

- PHP for interacting with the internal user and searches DataBases

- independent of anything

# Deployed version

❖ one of the deployed version

❖ Demo time
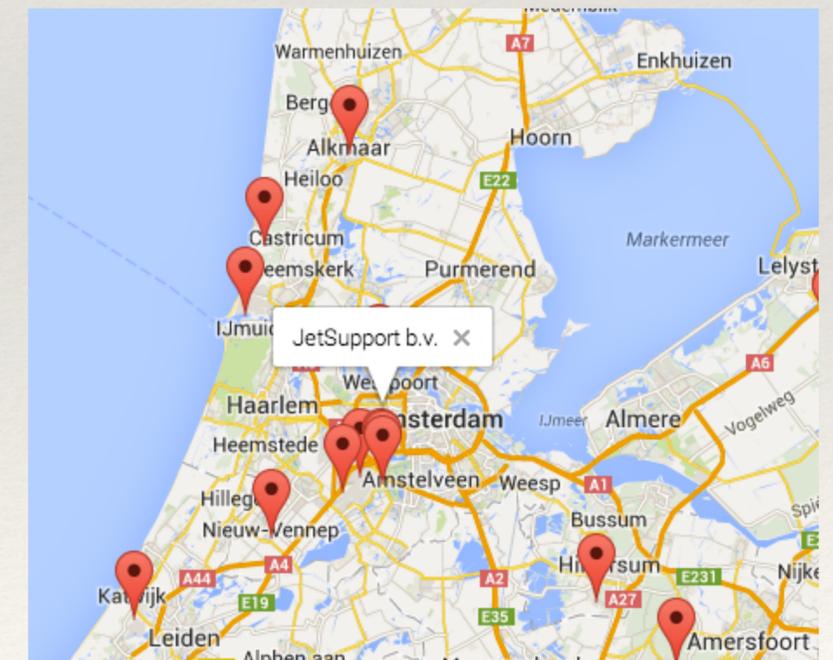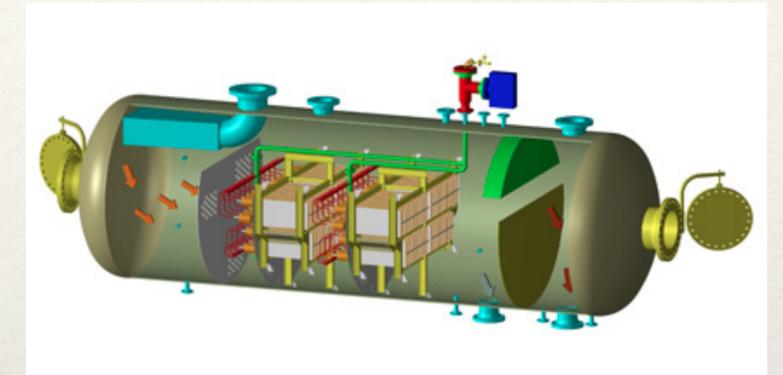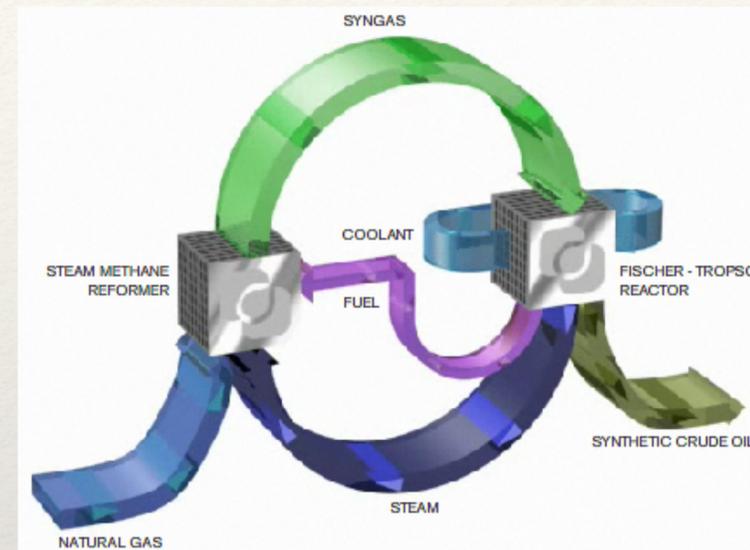
# CoDAK language

- ❖ Expanded keyword search (use Sesame to construct Solr query)

- ❖ locatie*, alles*, etc.

- ❖ Ex.

    - ❖ gas+Limburg (9 docs)

    - ❖ gas locatie*limburg  (11)

- ❖ Easily extendable, based on keyword-query table, no parsing required

# Other version of CoDAK



- ❖ Heat and Energy -related

- ❖ Separation technology -related

- ❖ Map integration

- ❖ Plain version (sandbox)

# Research opportunities

* major:

    * CoDAK language

    * Deep search map integration

    * "tamper" with Google

    * cluster documents/enhance ontology

* minor:

    * Ontology visualisation/interaction

    * Analyze user queries to enhance the ontology

    * Apply in a field of your own