# Why and how is this a "related document"?: Semantics-based analysis of and navigation through heterogeneous text corpora

**Bettina Berendt & Daniel Trümper
(KU Leuven / HU Berlin)
Blaž Fortuna, Marko Grobelnik & Dunja
Mladenič (JSI Ljubljana)**
www.cs.kuleuven.be/~berendt

Porpoise

File   Export   Help

Similarity Dimensions
Textual
Current   0.02000

Porpoise v. 0.1, 2007-12-04

ICT Motivation:
Global+local interaction; beyond "similar documents"

# Application motivation: Beyond dedicated search engines

**1. News and blogs**

**2. Multilingual sources**
➔ **Good results in semi-automatic ontology learning based on simple machine translation**



(Lloyd et al., *Proc. CAAW* 2006;
Berendt et al., *Kommunikation, Partizipation und Wirkungen im Social Web* , 2008;
Berendt, Fortuna et al., in prep.)

# PASCAL motivation:
# Re-use Textgarden's bread&butter and advanced tools

- **Text to bag-of-words**

- **Ontogen**



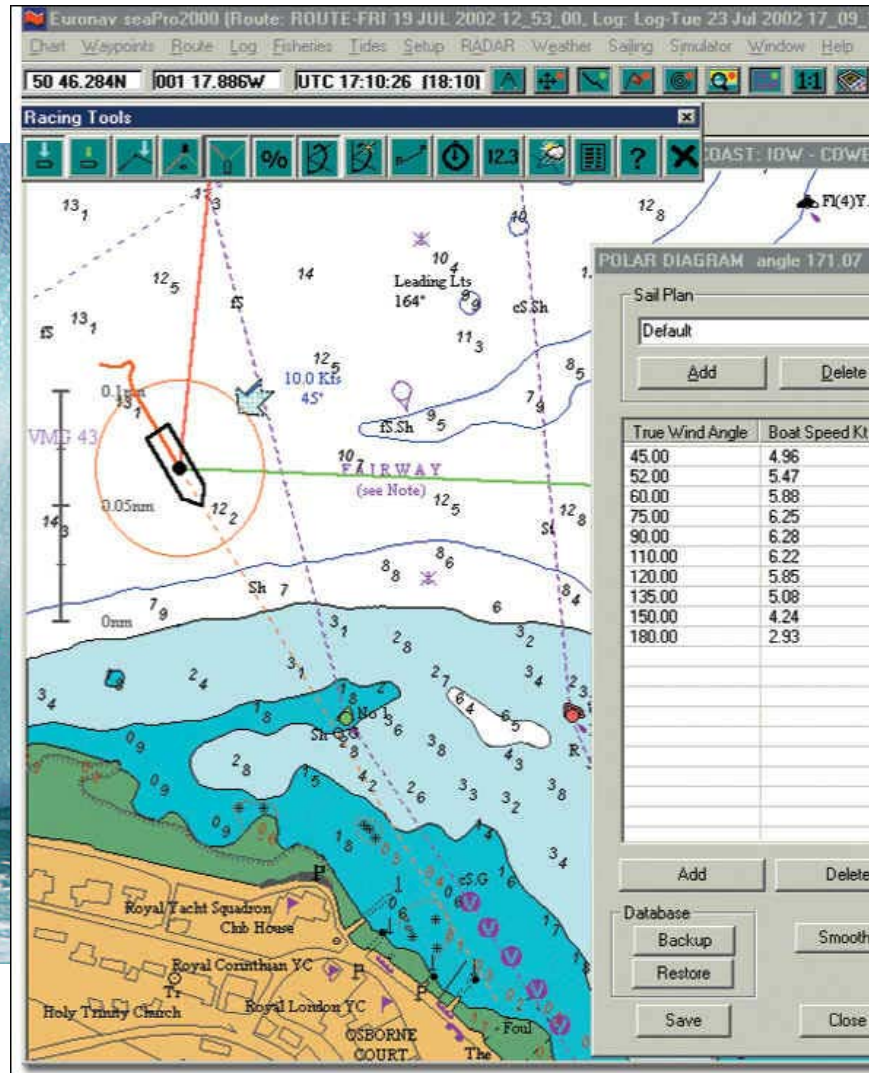http://www.textmining.net
http://ontogen.ijs.si/

# Solution vision:
## *PORPOISE – Sailing the Internet*

**Search**

**Global Analysis**

**Local analysis**

# Solution approach: Architecture & states overview

**Web**

**Search**

**Specify sources & filters ***

**Retrieval & Preprocessing ***

Data / tool
External
*Textgarden tool*
User action
Created in this project *

**Import ontology ***

**Global Analysis**

*Ont. Learning (Ontogen)*

**Source doc.s database***

*Build ontology*

**Select Document ***

**Select neighbour-hood ***

**Local analysis**

**Aspect-based similarity search***

**Refocus ***

**Construct composite-similarity neighbourhood ***

# Retrieval and preprocessing



- **Crawler / wrapper * (uses Blogdigger)**
- **Translator * (uses Babelfish)**
- **Preprocessing (*Txt2Bow*)**
- **NER (GATE)**
- **Similarity Computation ***

# Inspection of ontology and instances

# Inspection of documents

# More on documents

# The neighbourhood of a document

# Constructing the similarity measure & neighbourhood (II)

# Constructing the similarity measure & neighbourhood (III)

**Web**

Specify sources & filters *

**Retrieval & Preprocessing ***

Import ontology *

**Source doc.**

*Ont. Learning (Ontogen)*

*Build ontology*

Select Document *

Select neighbour-hood *

**Aspect-based similarity search***

Refocus *

Construct composite-similarity neighbourhood

**English blog**

**German blog**

**English news**

---

**Porpoise**

File   Export   Help

Lower Textual Similarity

Similarity Thresholds

Textual

Current   0.02400

Adjust threshold

Named Entities

Current   0.12000

Adjust threshold

Date

Current   0.5

Adjust threshold

Similarity 0.5 equals

4    days before and after the selected document

Published Before

Cluster   **Neighbours**

**A German-language blog**

**A news source**

**Most neighbours are English-language blogs**

# Comparing documents



**Porpoise**

File   Export   Help

Similarity Dimensions

Textual

---

**Document: Action: The Capitul**

| Name | feeds_dailykos_com_Action_The_Cap |
|---|---|
| Title | Action: The Capitulation Caucus |
| URL | http://feeds.dailykos.com/~r/dailykos |
| Type | BLOG |
| Language | English |
| Publication Date | Tue Sep 11 09:25:00 CEST 2007 |

**Content**   Original Content   Named Entities

to continue rubber stamping George Bu
And so it's time for the community to
dialing fingers and get to work.Where
Congressman stand? Check that. Where
Democratic member of Congress stand?
vote for H.R. 3087? Will they continu
majority of Democrats who want a time
now? Below the fold is the name and p
Democratic member of the House of Rep
goal here is to find
out, "Yes or no, does Congressman ___
Tell Congress NOT to support this cap
the fold, the name and phone number o
Representative, and if you are unsure
Representative is, information on how

---

**Document: 24th Chaos Communication Co...**

| Name | g_24th_Chaos_Communication_Congress_2007_Call_for_Participation.xml |
|---|---|
| Title | 24th Chaos Communication Congress 2007: Call for Participation |
| URL | .org/2007/24th-chaos-communication-congress-2007-call-for-participation/ |
| Type | BLOG_TRANS |
| Language | German |
| Publication Date | Fri Aug 24 03:49:00 CEST 2007 |

**Content**   Original Content   Named Entities

□ □ HomeAbout this blogImpressum □ "□" Friday 24 August
2007 24th
chaos Communication Congress 2007: Call for Participation
is hardly
the chaos Communication Camp past and already walks the
preparations
für the 24th chaos Communication Congress 2007 this like
each year of
27 30 December in Berlin to take place in front in
addition "Call
for a Participation" started Vorschläge für
program-guides now on
the Congress können by 12 October to be submitted in
addition gives

Close

ed Before

# Comparing documents; utilizing multilingual sources

**Porpoise**

File   Export   Help

Similarity Dimensions

Textual

---

**Document: Action: The Capitula...**

| | |
|---|---|
| Name | feeds_dailykos_com_Action_The_Cap |
| Title | Action: The Capitulation Caucus |
| URL | http://feeds.dailykos.com/~r/dailykos |
| Type | BLOG |
| Language | English |
| Publication Date | Tue Sep 11 09:25:00 CEST 2007 |

**Content** | Original Content | Named Entities

to continue rubber stamping George Bu
And so it's time for the community to
dialing fingers and get to work.Where
Congressman stand? Check that. Where
Democratic member of Congress stand?
vote for H.R. 3087? Will they continu
majority of Democrats who want a time
now? Below the fold is the name and p
Democratic member of the House of Rep
goal here is to find
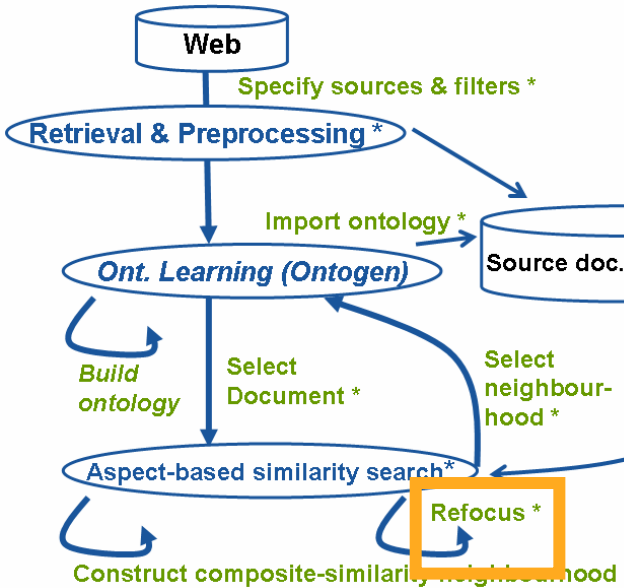out, "Yes or no, does Congressman ___
Tell Congress NOT to support this cap
the fold, the name and phone number o
Representative, and if you are unsure
Representative is, information on how

---

**Document: 24th Chaos Communication Co...**

| | |
|---|---|
| Name | g_24th_Chaos_Communication_Congress_2007_Call_for_Participation.xml |
| Title | 24th Chaos Communication Congress 2007: Call for Participation |
| URL | .org/2007/24th-chaos-communication-congress-2007-call-for-participation/ |
| Type | BLOG_TRANS |
| Language | German |
| Publication Date | Fri Aug 24 03:49:00 CEST 2007 |

Content | **Original Content** | Named Entities

HomeAbout this
blogImpressum

«

»

Freitag, 24. August 2007

24th Chaos
Communication Congress 2007: Call for Participation
Kaum ist
das Chaos Communication Camp vorbei, und schon schreiten
die Vorbereitungen für den 24th Chaos Communication
Congress 2007 voran. Dieser wird wie jedes Jahr vom
27.-30. Dezember in Berlin stattfinden. Dazu wurde jetzt
ein "Call for Participation" gestartet. Vorschläge für
Programmslots auf dem Congress können bis zum 12. Oktober

Close

ed Before

# Refocusing

# Structuring a neighbourhood



Web

Specify sources & filters *

Retrieval & Preprocessing *

Import ontology *

Ont. Learning (Ontogen)

Source doc

Build ontology

Select Document *

Select neighbour-hood *

Aspect-based similarity search*

Refocus *

Construct composite-similarity neighbourhood

File   Export   Help

Similarity Dimensions

Textual
Current   0.02000
Change similarity
-0.1   +0.1
-0.01   +0.01

Named Entities
Current   0.1200
Change similarity
-0.1   +0.1
-0.01   +0.01

Date
Current   0.02000
Change similarity
-0.1   +0.1
-0.01   +0.01

Lower Textual Similarity

Published Before

Cluster   Neighbours

# Ex.: Finding a "story"



**Document: Arizona congressman to retire ...**

**Porpoise**

File  Export  Help

Similarity Thresholds

Textual
Current  0.1000
Adjust threshold

Named Entities
Current  0.02000
Adjust threshold

Date
Current  0.5
Adjust threshold

Similarity 0.5 equals
4    days before and
     after the selected
     document

Lower Textual Similarity

Published Before

Published After

**Evaluation?
User studies!**

Cluster  **Neighbours**

Porpoise v. 0.1, 2008-01-20

# "Pump-priming": PORPOISE as catalyst

# Finally ... could I express it better?