**Slide 1**

**Sailing the Corpus Sea: Tools for Visual Discovery of Stories in Blogs and News**

**Bettina Berendt**

**www.cs.kuleuven.be/ ~berendt**



---

**Slide 2**

**About me**

Bettina Berendt [ Edit ]
Associate Professor at Katholieke Universiteit Leuven [ Edit ]
Brussels Area, Belgium [ Edit ]

Profile | Q&A | Recommendations | Conne...

Current • **Associate Professor at** Leuven [ Edit ]

Past • Assistant Professor at H...
• Research positions (Deta...
Humboldt Univ. Berlin, Ur...
Hamburg

Education • Humboldt-Universität zu B...
• Universität Hamburg
• The University of Edinburg...
• Freie Universität Berlin
• University of Cambridge

**Data Mining for Enlightenment**

Bettina Berendt
www.cs.kuleuven.be/ ~berendt

---

**Slide 3**

**Thanks to ...**

**Daniel Trümper**

**Tool at http://www.cs.kuleuven.be/~berendt/PORPOISE/**

**Ilija Subašić**

**Tool forthcoming;**

**all beta testers and experiment participants welcome!**

---

**Slide 4**

**First motivation:**
**Global+local interaction; beyond "similar documents"**

**with respect to what?**

---

**Slide 5**

**Solution vision:**
*Sailing the Internet*

**Search**

**Global Analysis**

**Local analysis**

---

**Slide 6**

**Solution approach: Architecture & states overview (version 1)**

**Web**

**Search**          **Specify sources & filters**

**Retrieval & Preprocessing**

**Import ontology**

**Global Analysis**

➔ **Story space**

**Ontology Learning**          **Source doc.s database**

*Build ontology*          **Select Document**          **Select neighbour-hood**

**Local analysis**

➔ **Document space**

**Aspect-based similarity search**          **Refocus**

**Construct composite-similarity neighbourhood**

---

## Retrieval and preprocessing



- Crawler / wrapper (uses Yahoo! / Google News; Blogdigger)
- Translator (uses Babelfish)
- Preprocessing (uses Textgarden, Terrier)
- Named-entity recognition (uses GATE, OpenCalais)
- Similarity Computation

## Ontology learning (1)



Tool: Blaž Fortuna: http://blazfortuna.com/projects/ontogen

## Ontology learning (2)

## Inspection of ontology and instances

## Inspection of documents

## More on documents

The neighbourhood of a document



Constructing the similarity measure & neighbourhood (I)



Constructing the similarity measure & neighbourhood (II)



Constructing the similarity measure & neighbourhood (III)



Comparing documents



Comparing documents; utilizing multilingual sources

**19**

## Refocusing

Porpoise
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

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

---

**20**

## Structuring a neighbourhood

Porpoise
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

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

---

**21**

## Ex.: Finding a "story"

Document: Arizona congressman to retire …

Porpoise
File  Export  Help

Similarity Thresholds
Textual
Current  0.1000
Adjust threshold
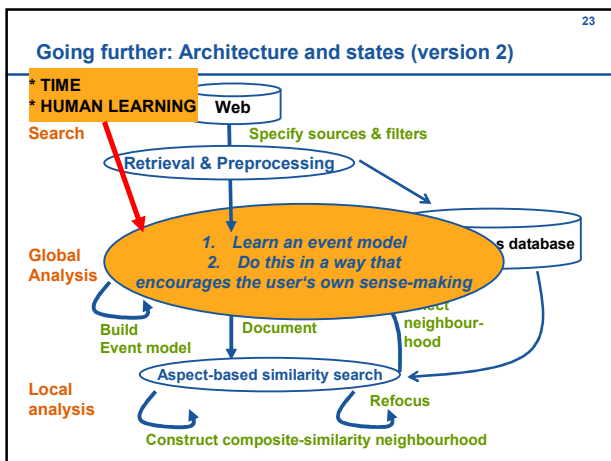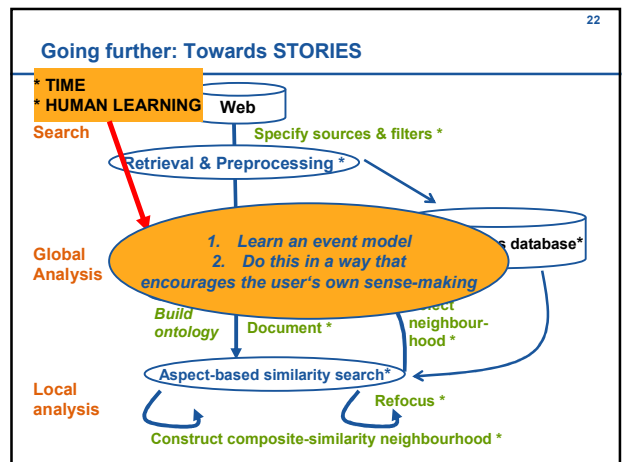
Named Entities
Current  0.02000
Adjust threshold

Date
Current  0.5
Adjust threshold

Similarity 0.5 equals
days before and after the selected document

Lower Textual Similarity

Published Before          Published After

Cluster  Neighbours

Porpoise v. 0.1, 2008-01-20

---

**22**

## Going further: Towards STORIES

* TIME
* HUMAN LEARNING
Web

Search
Specify sources & filters *
Retrieval & Preprocessing *

Global Analysis
1. Learn an event model
2. Do this in a way that encourages the user's own sense-making
s database*

Build ontology
Document *
Select neighbour-hood *

Local analysis
Aspect-based similarity search*
Refocus *
Construct composite-similarity neighbourhood *

---

**23**

## Going further: Architecture and states (version 2)

* TIME
* HUMAN LEARNING
Web

Search
Specify sources & filters
Retrieval & Preprocessing

Global Analysis
1. Learn an event model
2. Do this in a way that encourages the user's own sense-making
s database

Build Event model
Document
Select neighbour-hood

Local analysis
Aspect-based similarity search
Refocus
Construct composite-similarity neighbourhood

---

**24**

## Solution approach 1: Find latent topics

Document Atlas 2.0 -- Text Garden
File  About

Map Properties
Document names
Common words
Magnifying glass
Gradient
Font sizes:

Categories
Category:  august

Thresholds
Documents:
Relations:

Selected Document
Title:
august/2007-08-04- id 329
Content:
madeleine sight the Daily Express and the Daily Mirror both report a claim that miss four-year-old Madeleine_McCann be see in Belgium . the Express describe the development as one of the first meaningful event since the abduction . it must give

ROBERT_MURAT, SUSPECT, POLICE, REPORT, SEARCH, NIGHT, MAN, QUESTION, APARTMENT, FRIEND

• temporal development only by comparative statics
• no „drill down" possible
• no fine-grained relational information
➔ lacks structure

Tool: Blaž Fortuna : http://docatlas.ijs.si

## Solution approach 2: Temporal latent topics



Figure 6: Theme evolution graph for Asia Tsunami

- no fine-grained relational information
- "themes" are fixed by the algorithm
- no „drill down" possible
➔ no combination of machine and human intelligence

Mei & Zhai, *PKDD 2005*

---

## The ETP3 problem

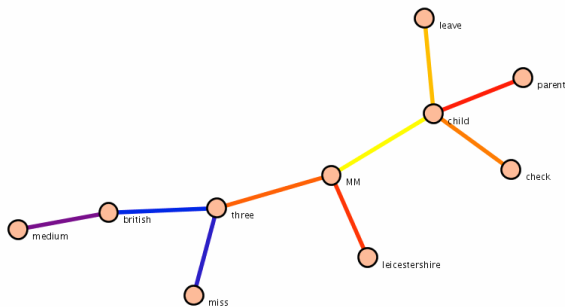Evolutionary theme patterns discovery, summary & exploration

1. identify topical sub-structure in a set (generally, a time-indexed stream) of documents constrained by being about a common topic
2. show how these substructures emerge, change, and disappear (and maybe re-appear) over time
3. give users intuitive and interactive interfaces for exploring the topic landscape and the underlying documents

   *and for their own sense-making*

   *– use machine-generated summarization only as a starting point!*

---

## Ingredients of a solution
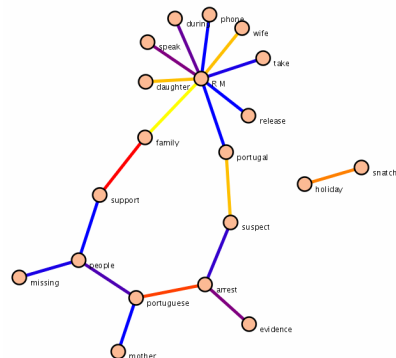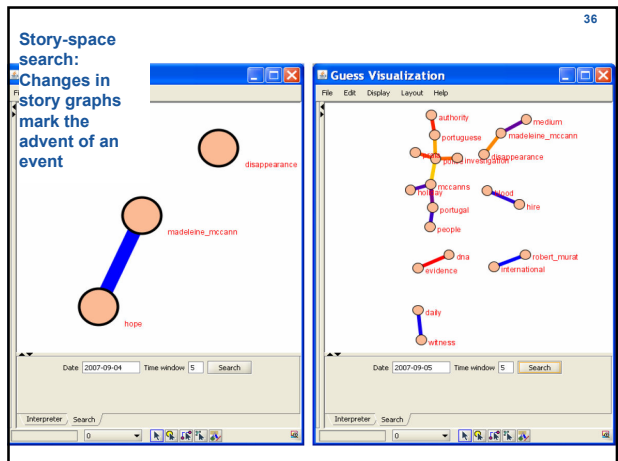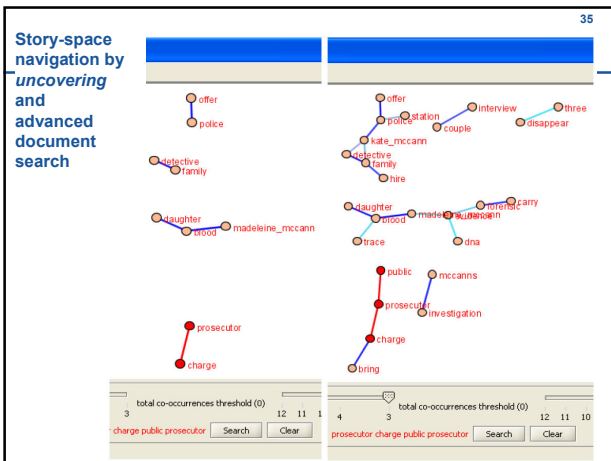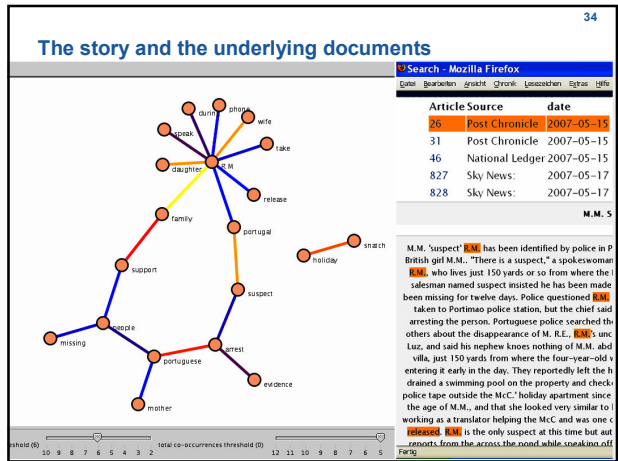**to ETP3 =** Evolutionary theme patterns discovery, summary and exploration

**Document / text pre-processing**
- **Template recognition**
- **Multi-document named entities**
- **Stopword removal, lemmatization**

**Document summarization strategy**
- **no topics, but salient concepts & relations**
- **time window; word-span window**

**Selection approach for concepts**
- **concepts = words or named entities**
- **salient concept = high TF & involved in a salient relation, time-indexed**

**Similarity measure to determine relations**
- **bursty co-occurrence**

**Burstiness measure**
- **time relevance, a "temporal co-occurrence lift"**

$$TR_i(b_1,b_2) = \frac{freq_i(b_1,b_2)}{freq_T(b_1,b_2)}.$$

**Interaction approach**
- **Graphs (& layout)**
- **Comparative statics or morphing**
- **Drill-down: "uncovering" relations**
- **Links to documents (in progress)**

STORIES

---

## "Powerpoint demo"

---

## An event: a missing child



---

## A central figure emerges in the police investigations

**Uncovering more details**

**Uncovering more details**

**An eventless time**

**The story and the underlying documents**

**Story-space navigation by *uncovering* and advanced document search**

**Story-space search: Changes in story graphs mark the advent of an event**

## Data collection and preprocessing

- Articles from Google News 05/2007 – 11/2007 for search term "madeleine mccann"
  - *(there was a Google problem in the December archive)*
- Only English-language articles
- For each month, the first 100 hits
- Of these, all that were freely available ➜ 477 documents

- Preprocessing:
  - HTML cleaning
  - tokenization
  - stopword removal

## Story elements

- *content-bearing words*
  - the 150 top-TF words without stopwords

$$TF_{ij} = \frac{n_{ij}}{|d_i|}$$

## Story stages: co-occurrence in a window

[Sep 7, 2007] Madeleine McCann's mother is being made a **formal suspect** by police investigating the child's disappearance. The move was revealed as Kate McCann, 39, who emerged early on Friday morning after almost 11 hours of questioning by Portuguese police on Thursday night, prepared to face officers again.

**"mother" and "suspect" co-occur**
- **in a window of size ≥ 6 (all words)**
- **in a window of size ≥ 2 (non-stopwords only)**
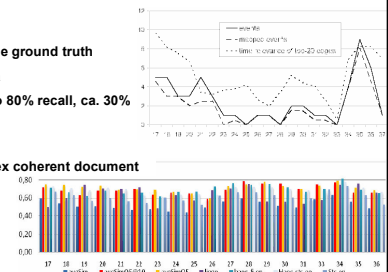
## Salient story elements

1. Split whole corpus T by week (17 = 30 Apr + until 44 = 12 Nov +)
2. For each week
   - Compute the weights for corpus t for this week
3. Weight =
   - *Support* of co-occurrence of 2 content-bearing words $w_1$, $w_2$ in t =
   (# articles from t containing both $w_1$, $w_2$ in window) / (# all articles in t)
4. Threshold
   - Number of occurrences of *co-occurrence($w_1$, $w_2$)* in t ≥ $\theta_1$ (e.g., 5)
   - *Time-relevance* TR of *co-occurrence($w_1$, $w_2$)* =
   support(*co-occurrence($w_1$, $w_2$)*) in t / support(*co-occurrence($w_1$, $w_2$)*) in T ≥ $\theta_2$ (e.g., 2) *
5. Rank by TR, for each week identify top 2
6. Story elements = peak words = all elements of these top 2 pairs (# = 38)

## Salient story stages, and story evolution

7. Story stage = co-occurrences of peak words in t
   - For each week t: aggregate over t-2, t-1, t ➜ moving average

8. Story evolution = how story stages evolve over the t in T

## Evaluations *(so far ...)*

1. Information retrieval quality
   - Challenge: What is the ground truth
   - ➜ Build on Wikipedia
   - Edges – events: up to 80% recall, ca. 30% precision
2. Search quality
   - Story subgraphs index coherent document clusters

3. Learning effectiveness
   - Document search with story graphs leads to averages of
   - 75% accuracy on judgments of story fact truth
   - 3.4 nodes/words per query

7

## Summary

Navigation in story space → story building

+

Document search

+

Navigation in document space

**lead to understandable, useful + intuitive interfaces**

## The Future

- **Better language processing**
- **Linkage information!**
- **Opinion mining**

# Thanks!

## References

Subašić, I. & Berendt, B. (2008). Web Mining for Understanding Stories through Graph Visualisation. In *Proceedings of ICDM 2008*. IEEE Press.

Berendt, B. and D. Trümper (in press). Semantics-based analysis and navigation of heterogeneous text corpora: the Porpoise news and blogs engine. In I.-H. Ting & H.-J. Wu (Eds.), *Web Mining Applications in E-commerce and E-services*, Berlin etc.: Springer.

Berendt, B. & Subašić, I. (in press). Measuring graph topology for interactive temporal event detection. To appear in *Künstliche Intelligenz*.

Berendt, B. & Subašić, I. (under review). Discovery of interactive graphs for understanding and searching time-indexed corpora.

Please see http://www.cs.kuleuven.be/~berendt/ for these papers.