# Cross-Language Ideology Detection Using Linguistics-Based Knowledge Extraction

## Department 06641
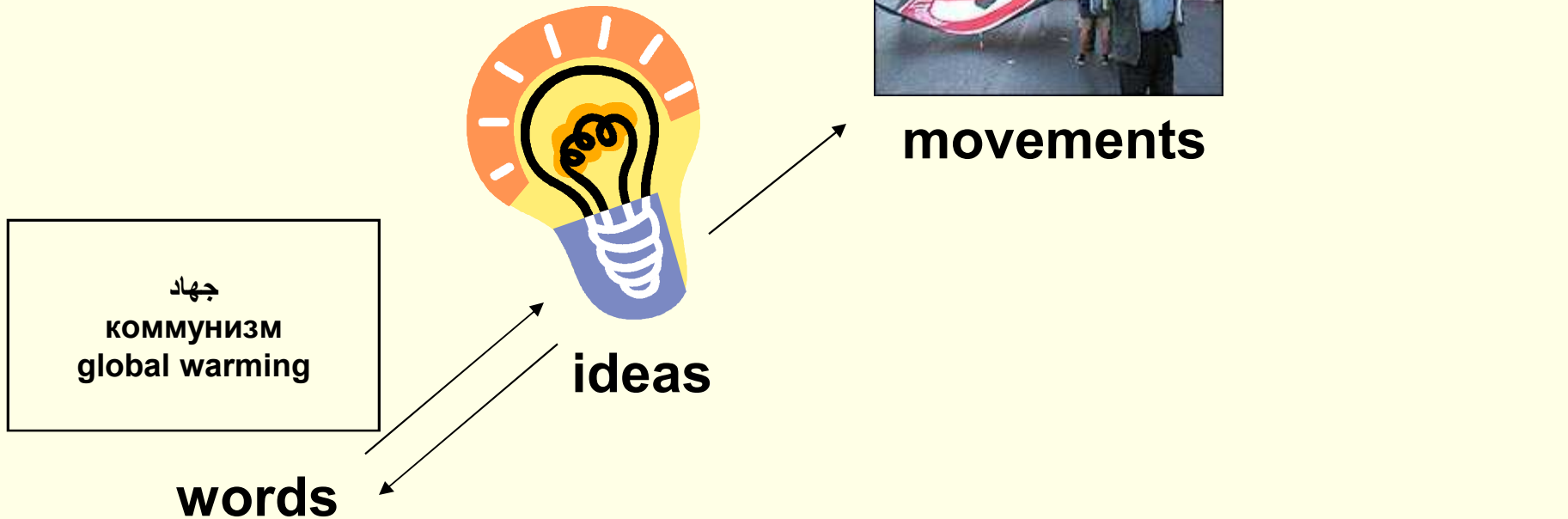
### September 13, 2006

**Peter Chew, Stephen Verzi, JT McClain, Travis Bauer**
**Sandia National Laboratories**

Sandia National Laboratories

# Background

Words, or language, persuade people to join movements which can in turn become threats.

**threats**

**movements**

جهاد
коммунизм
global warming

**words**

**ideas**

Sandia National Laboratories

# Our approach

- **Representative sample of web pages, with focus towards blogs, commentaries**
  - **Multi-lingual (English, Russian, Arabic, French, Spanish)**
- **Collect in a SQL Server database**
- **Analyze, based on language and/or links, the prominence of different clusters**
  - **We expect approach to be 'blind' to language**
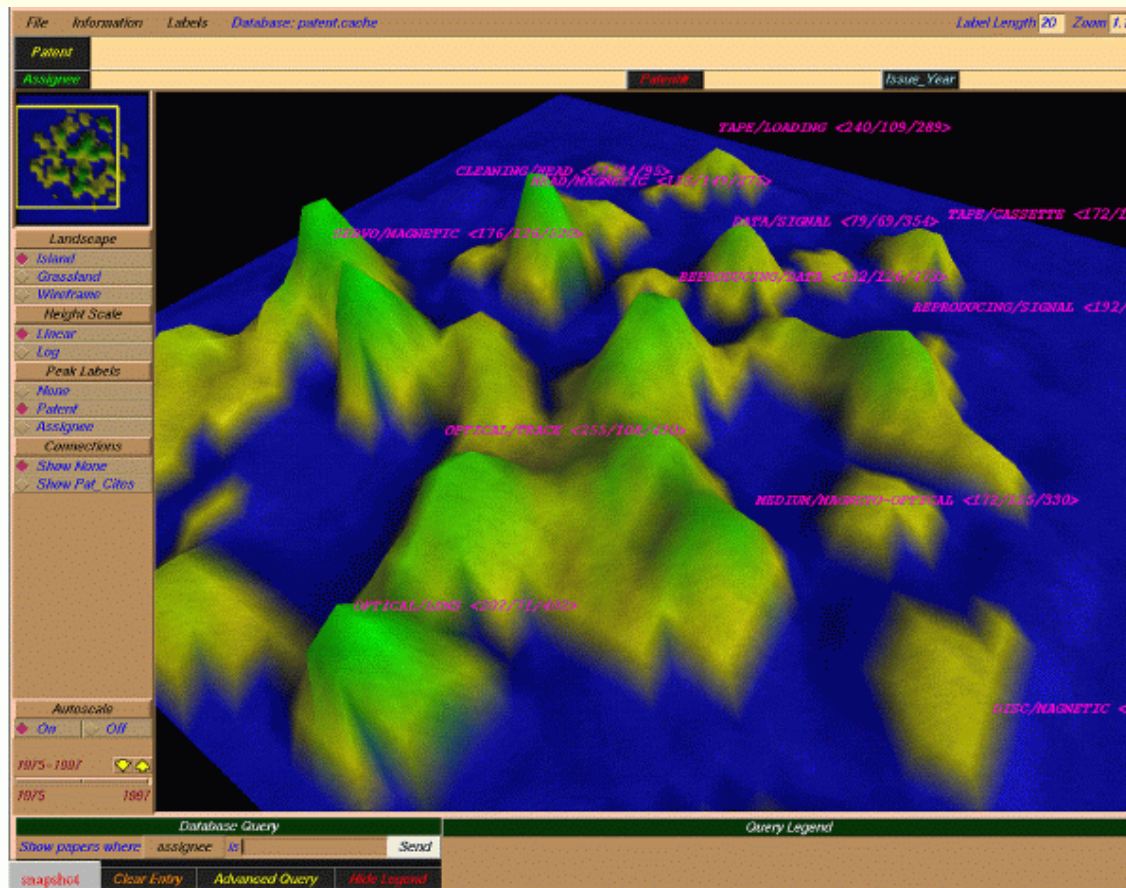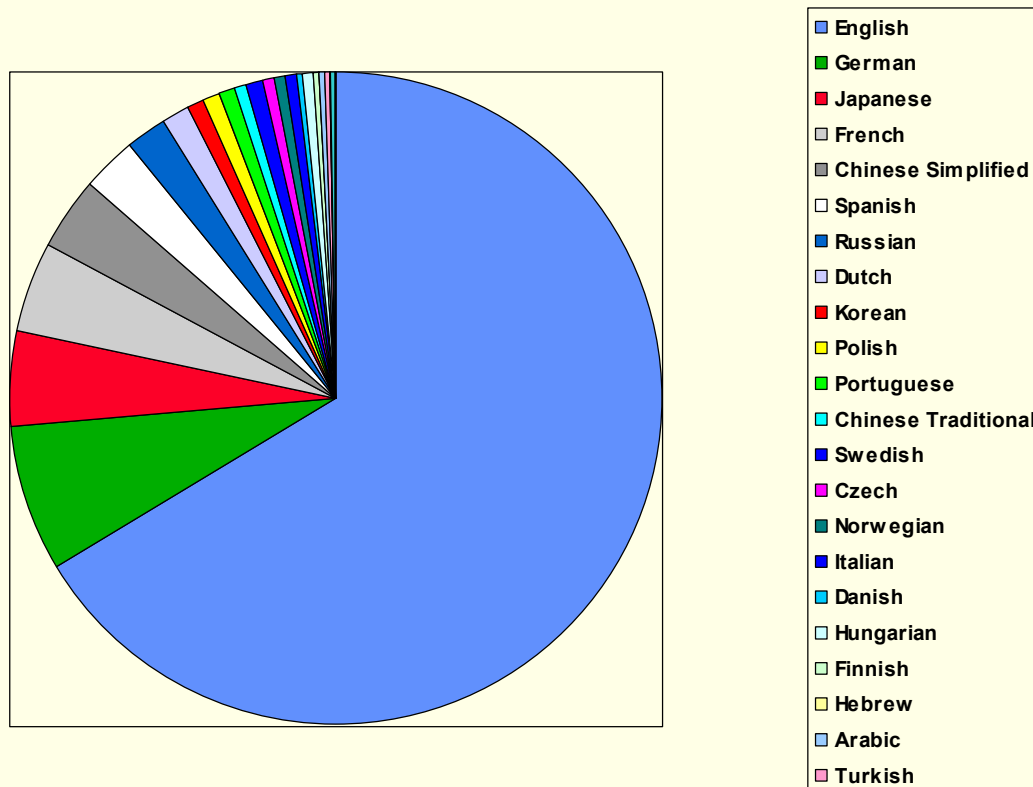- **Analyze changes over time**

# What we are aiming for – an example



Figure 2: The mountain terrain metaphor provides an intuitive exploration environment.

# Distribution of internet content by language



Legend:
- English
- German
- Japanese
- French
- Chinese Simplified
- Spanish
- Russian
- Dutch
- Korean
- Polish
- Portuguese
- Chinese Traditional
- Swedish
- Czech
- Norwegian
- Italian
- Danish
- Hungarian
- Finnish
- Hebrew
- Arabic
- Turkish

# Data collection to date

| | Number of documents | % in sample | approx. % in population |
|---|---|---|---|
| **English** | 176,307 | **50.7** | **66.5** |
| **Russian** | 62,954 | **18.1** | **2.0** |
| **Arabic** | 30,248 | **8.7** | **0.2** |
| **Spanish** | 39,898 | **11.5** | **2.6** |
| **French** NEW! | 38,479 | **11.1** | **4.4** |
| *Others* | - | - | **24.3** |
| **TOTAL** | **347,886** | **100.0** | **100.0** |

# Approaches to cross-language information retrieval

- **Translate the query**
  - **Efficacy is constrained by quality of machine translation**



- **Train algorithm on parallel corpora**
  - **Translations should:**
    - **Be available in target languages**
    - **Be reliable**
    - **Be sufficiently large in size**
    - **Cover target subject domain**
    - **Be free of undue copyright restrictions**
    - **Be electronically available**
    - **Be alignable**

# The Bible as a Parallel Corpus

- **Resnik, Olsen & Diab (1999) showed that the Bible fulfills all of these criteria and is surprisingly suitable as a parallel corpus**
  - **Translations in > 2,400 languages and rising**
  - **Great care taken over translations**
  - **Respectably large compared to other corpora**
  - **Covers many modern genres**
  - **Covers up to 85% of modern vocabulary**
  - **Generally free of copyright restrictions**
  - **Electronically available**
  - **Alignable**

# Language coverage - detail

## A STATISTICAL SUMMARY OF LANGUAGES WITH THE SCRIPTURES

A summary, by geographical area and type of publication, of the number of different languages and dialects in which publication of at least one book of the Bible had been registered as of December 31, 2005.

| Continent/Region | Portions | Testaments | Bibles | Bibles, DC* | Total |
|---|---|---|---|---|---|
| Africa | 223 | 301 | 159 | (29) | 683 |
| Asia | 218 | 244 | 131 | (28) | 593 |
| Australia/New Zealand/ Pacific Islands | 148 | 234 | 38 | (9) | 420 |
| Europe | 114 | 36 | 61 | (47) | 211 |
| North America | 39 | 30 | 7 | (0) | 76 |
| Caribbean Islands / Central America / Mexico/South America | 118 | 270 | 29 | (9) | 417 |
| Constructed Languages | 2 | 0 | 1 | (0) | 3 |
| TOTALS | 862 | 1,115 | 426 | (122) | 2,403 |

* This column is a sub-section of the Bibles column – for example, there is a translation of the Deuterocanon for 47 of the 61 languages of Europe in which the Bible has been translated.

[A few corrections were made to our language databases and are reflected in this statistical summary]

Per http://www.biblesociety.org/latestnews/latest341-slr2005stats.html

# The 'Unbound Bible'



85 translations (some partial) in 51 languages, in common format

# The 'Unbound Bible' – a sample



```
web_utf8.txt - Notepad
File  Edit  Format  View  Help
43N          1          1          In the beginning was the Word, and
43N          1          2          The same was in the beginning with
43N          1          3          All things were made through him.
43N          1          4          In him was life, and the life was
43N          1          5          The light shines in the darkness,
43N          1          6          There came a man, sent from God, w
43N          1          7          The same came as a witness, that h
43N          1          8          He was not the light, but was sent
43N          1          9          The true light that enlightens eve
43N          1          10         He was in the world, and the world
43N          1          11         He came to his own, and those who
43N          1          12         But as many as received him, to th
43N          1          13         who were born, not of blood, nor o
43N          1          14         The Word became flesh, and lived a
43N          1          15         John testified about him. He cried
43N          1          16         From his fullness we all received
43N          1          17         For the law was given through Mose
43N          1          18         No one has seen God at any time.
43N          1          19         This is John's testimony, when the
43N          1          20         He confessed, and didn't deny, but
43N          1          21         They asked him, "What then? Are yo
```
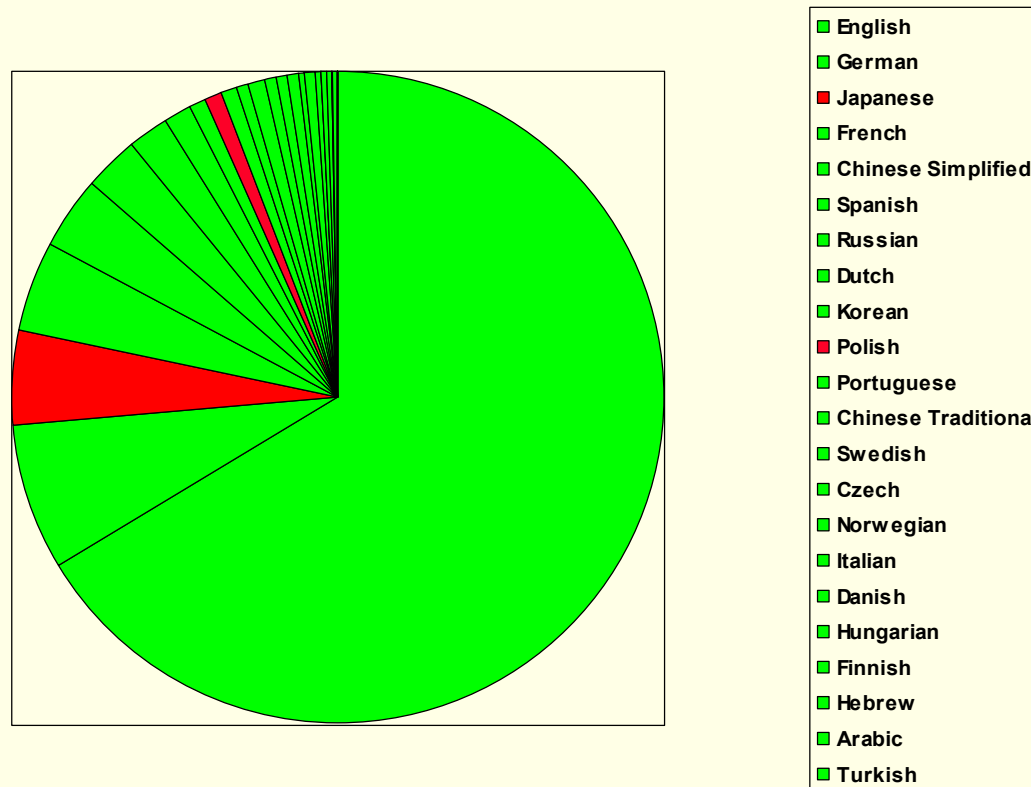
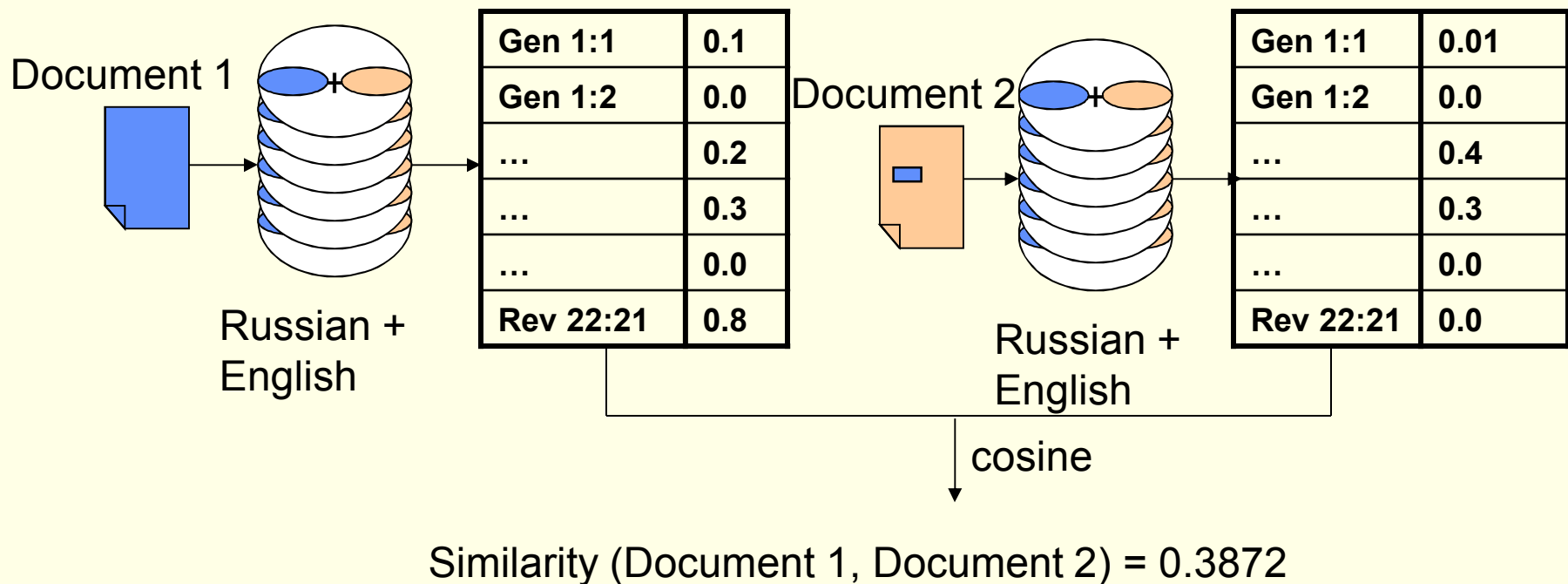# Coverage of internet content based on 'Unbound Bible'



- English
- German
- Japanese
- French
- Chinese Simplified
- Spanish
- Russian
- Dutch
- Korean
- Polish
- Portuguese
- Chinese Traditional
- Swedish
- Czech
- Norwegian
- Italian
- Danish
- Hungarian
- Finnish
- Hebrew
- Arabic
- Turkish

Sandia National Laboratories

**Cognitive Systems** Applying computer models of human cognition to create unique technology solutions

# Cross-Language Comparison – Method 1

- **Use separate STANLEY textual model for each language**

| | |
|---|---|
| **Gen 1:1** | 0.1 |
| **Gen 1:2** | 0.0 |
| … | 0.2 |
| … | 0.3 |
| … | 0.0 |
| **Rev 22:21** | 0.8 |

Document 1 → Russian

| | |
|---|---|
| **Gen 1:1** | 0.0 |
| **Gen 1:2** | 0.0 |
| … | 0.4 |
| … | 0.3 |
| … | 0.0 |
| **Rev 22:21** | 0.0 |

Document 2 → English

cosine

Similarity (Document 1, Document 2) = 0.3850

Sandia National Laboratories

# Cross-Language Comparison – Method 2

- **Use single STANLEY textual model of concatenated language document chunks**

Document 1

| Gen 1:1 | 0.1 |
|---------|-----|
| Gen 1:2 | 0.0 |
| … | 0.2 |
| … | 0.3 |
| … | 0.0 |
| Rev 22:21 | 0.8 |

Russian + English

Document 2

| Gen 1:1 | 0.01 |
|---------|------|
| Gen 1:2 | 0.0 |
| … | 0.4 |
| … | 0.3 |
| … | 0.0 |
| Rev 22:21 | 0.0 |

Russian + English

cosine

Similarity (Document 1, Document 2) = 0.3872

# Validation results: Test set in training set

**Method:** Index on entire Bible, measure average uninterpolated precision at doc. 1 for 66 books of Bible

|  |  | To | | | | |
|---|---|---|---|---|---|---|
|  |  | **Arabic** | **English** | **French** | **Russian** | **Spanish** |
| **From** | **Arabic** | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 |
|  | **English** | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 |
|  | **French** | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 |
|  | **Russian** | 1.00 | .99 | 1.00 | 1.00 | 1.00 |
|  | **Spanish** | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 |

# Validation results: Test set not in training set

**Method:** Index on entire Bible, obtain matrix of similarity measures for 5 conference abstracts where English and Spanish translations exist

|  |  | English | | | | |
|---|---|---|---|---|---|---|
|  |  | **Doc 1** | **Doc 2** | **Doc 3** | **Doc 4** | **Doc 5** |
| **Spanish** | **Doc 1** | **.607** | .043 | .045 | .035 | .022 |
|  | **Doc 2** | .049 | **.397** | .038 | .082 | .166 |
|  | **Doc 3** | .030 | .050 | .045 | .101 | .049 |
|  | **Doc 4** | .102 | .096 | .080 | **.189** | .105 |
|  | **Doc 5** | .035 | .131 | .039 | .042 | **.168** |

We have also used the framework successfully for Maori, to distinguish between the Treaty of Waitangi and the New Zealand National Anthem

Sandia National Laboratories

# Validation results: Test set not in training set

**Method:** Index on entire Bible, measure Mean Average Precision for 114 suras of Quran in English, Arabic, Russian, and Spanish (results comparable to McNamee & Mayfield 2004)

.35

(using 5-grams)

| - LSA<br>- no removal of stopwords | | To | | | |
| --- | --- | --- | --- | --- | --- |
| | | Arabic | English | Russian | Spanish |
| **From** | **Arabic** | 1.00 | .71  .60 | .62  .33 | .72  .46 |
| | **English** | .71  .49 | 1.00 | .90  .75 | .90  .53 |
| | **Russian** | .56  .40 | .92  .68 | 1.00 | .67  .45 |
| | **Spanish** | .66  .46 | .87  .78 | .74  .62 | 1.00 |

Method 1: Separate index for each language

Method 2: Single index for all languages

Sandia National Laboratories

# Precision-recall graphs for different test parameters

# Observations

- **Cross-language approach is easily extensible to new languages and corpora, including minority languages**

- **Resources exist which allow large parallel corpora to be built up from scratch in hours, at no monetary cost**

- **Unsurprisingly, the larger the training set, the better the precision-recall results**

- **Results appear to be comparable to, or better than, those achieved by other methods and reported recently in research literature; further testing may be needed**

- **Our best results were obtained using LSA, a single index for all languages, and without removing stopwords. This has the advantage of requiring no language-specific expertise to set up.**

Sandia National Laboratories

# Beyond the 'bundle of words'

He    wants    to    be    happy

**The 'bundle of words'**

# Of course, order does matter in language

# Why does the bundle of words approach work at all?

**Assume:**

– **vocabulary of 100,000 ($10^5$) words**

– **sentences are 20 words long**

– **word order important only within sentences**

$\Rightarrow$ **Contribution (in bits) to passage 'information':**

| | | | |
|---|---|---|---|
| – **From word *choice*:** $\log_2(20!)$ | $\approx$ | **61·0774** | **84.47%** |
| – **From word *order*:** $\log_2((10^5)^{20})$ | $\approx$ | **332·1928** | **15.53%** |
| | | | |
| – **Total** $\log_2(20! \times 10^{100})$ | $\approx$ | **393.2702** | **100.00%** |

# N-Word-Grams

- **Consecutive N-Word Lexical Entities**
- **Statistical Windowing Technique**
  - **window width**
  - **step size**

- **We are at Sandia National Laboratories …**

# Document-to-Document Similarities

• **All 120 permutations (n = 1 word-grams)**

# Document-to-Document Similarities

- **All 120 permutations (n <= 2 word-grams)**

# Document-to-Document Similarities

- **All 120 permutations (n = 2 word-grams)**

# Document-to-Document Similarities

- **All 120 permutations (2 <= n <= 3 word-grams)**

# Document-to-Document Similarities

- **All 120 permutations (n <= 5 word-grams)**

# Document-to-Document Similarities

- **All 120 permutations (n = 5 word-grams)**

# Ideological documents of known threat level

# Ideological documents of known threat level

# Detecting Ideology

Words, or language, persuade people to join movements which can in turn become threats.



**threats**



**movements**



**ideas**

جهاد
**коммунизм
global warming**

**words**

Sandia National Laboratories

# Similar Web Documents
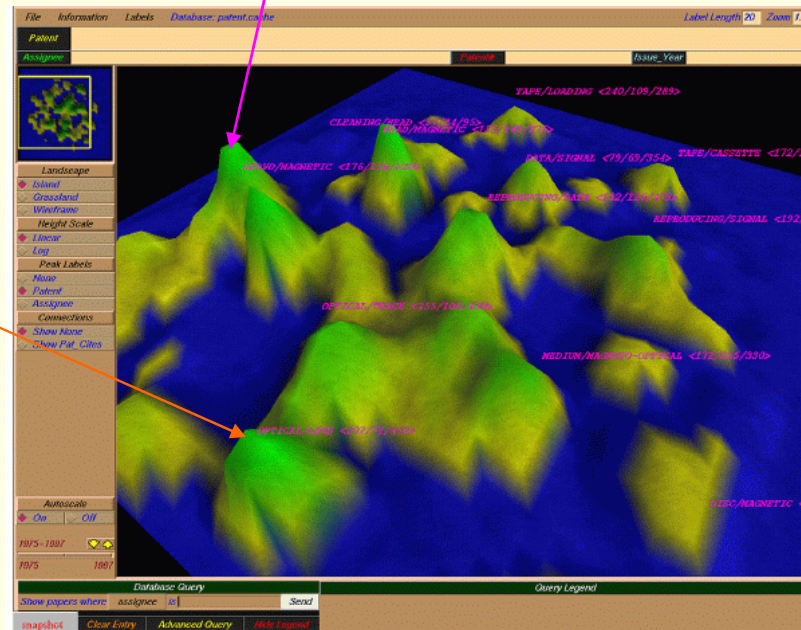
*Communist Manifesto*

*Mein Kampf*



Figure 2: The mountain terrain metaphor provides an intuitive exploration environment.

# Future directions

- **Continue to test CLIR algorithm to identify its strengths and weaknesses**
- **What chunk size yields best cross-language IR results, and why?**
- **Can we use the output of cross-language comparison to characterize documents by their ideology?**
- **Add further languages for more coverage of WWW**
- **Linguistic analysis beyond the 'bundle of words' approach**

# QUESTIONS?

Sandia National Laboratories

# Selected References

- Ackland, R. 2005. Mapping the U.S. Political Blogosphere: Are Conservative Bloggers More Prominent? mimeo., The Australian National University.

- Chomsky, Noam. 1988 *Language and Politics*. C.P. Otero (ed.). Montreal: Black Rose Books.

- Landauer, Thomas. 1998. An Introduction to Latent Semantic Analysis. In *Discourse Processes* 25, 259-284.

- McNamee, Paul, and James Mayfield. 2004. Character *N*-Gram Tokenization for European Language Text Retrieval. *Information Retrieval* 7:73-97.

- Resnik, Philip, Mari Olsen, and Mona Diab. 1999. The Bible as a Parallel Corpus: Annotating the Book of 2000 Tongues. In *Computers and the Humanities* 33, 1-2. pp 129-153.

- Resnik, Philip, and Noah Smith. 2003. The Web as a Parallel Corpus. In *Computational Linguistics* 29, No. 3, pp. 349-380.

# Selected References – "Bundle of Words"

- Thomas K. Laudauer, Darrell Laham, Bob Rehder, and M. E. Schreiner, "How Well Can Passage Meaning be Derived without Using Word-Order? A Comparison of Latent Semantic Analysis and Humans," Proceedings of the 19[th] Annual Meeting of the Cognitive Science Society, 1997.

- Fernando Pereria, "Formal Grammar and Information Theory: Together Again?," Philosophical Transactions of the Royal Society Series A 358, 1239-1253, 2000.

- Lynette Hirschman, Marc Light, Eric Breck, and John D. Burger, "Deep Read: A Reading Comprehension System," Proceedings of the ACL, 1999.
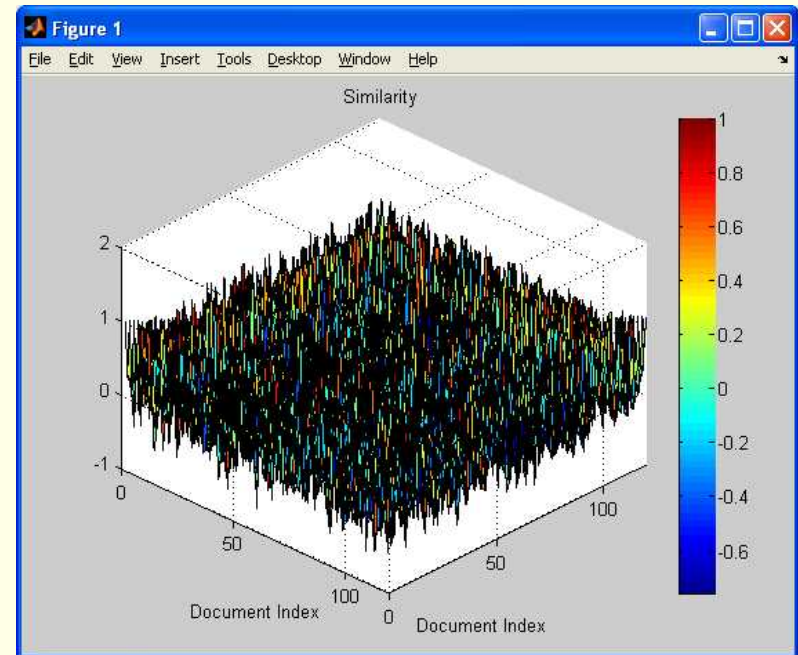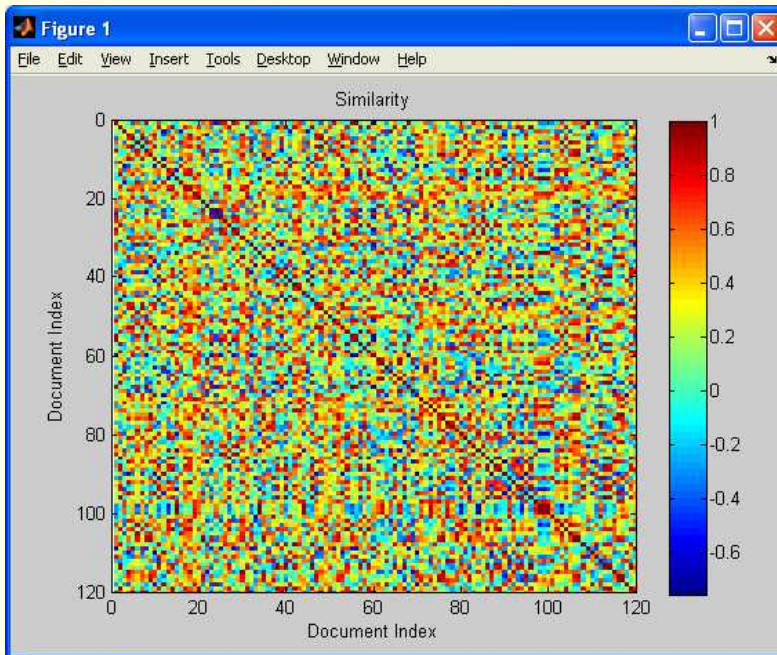
# Selected References – "Word-Grams"

- Roger Zhang, "A Simple C++ N-gram Extraction Package", Dalhousie University, Nova Scotia, Canada, November 28, 2004.

# Backup Slides

Sandia National Laboratories

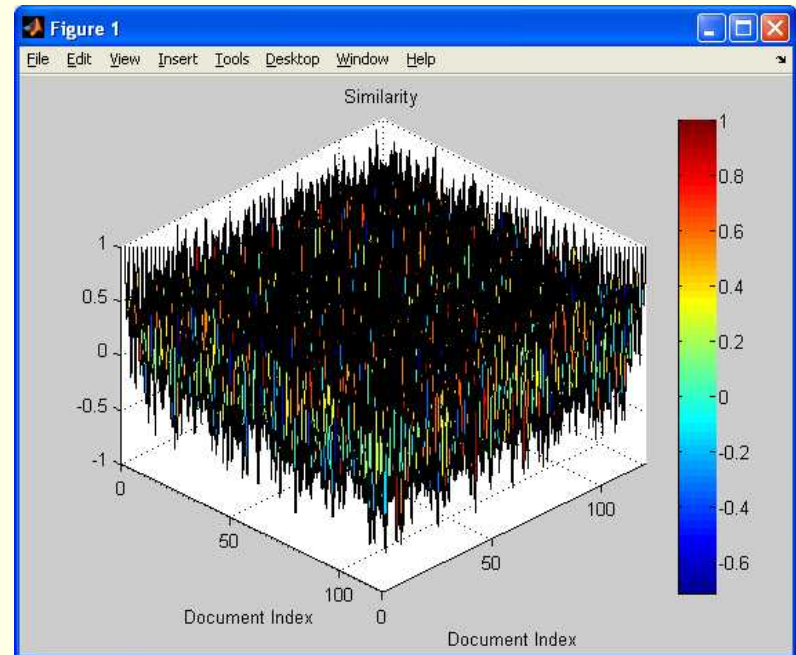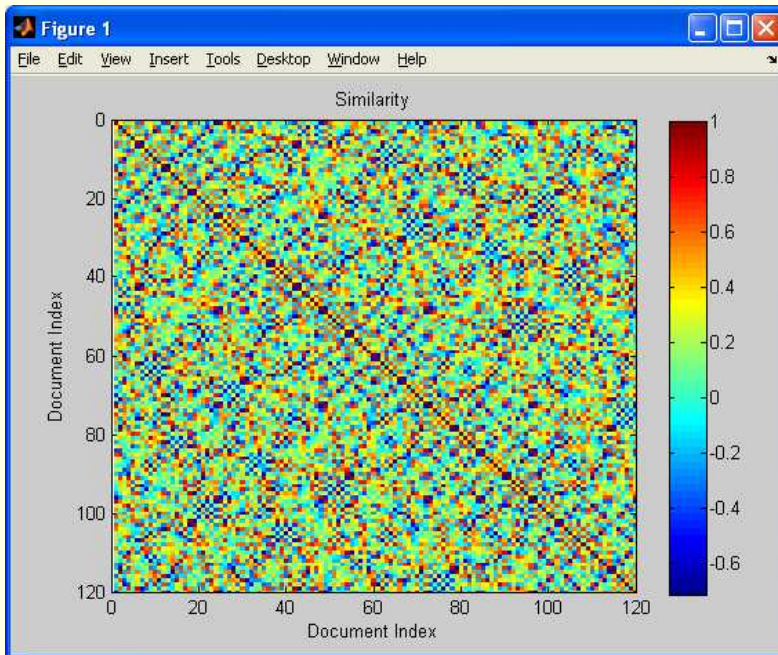# Document-to-Document Similarities

- **All permutations (n = 1 word-grams) - LSA**

# Document-to-Document Similarities

• **All permutations (n <= 2 word-grams) - LSA**
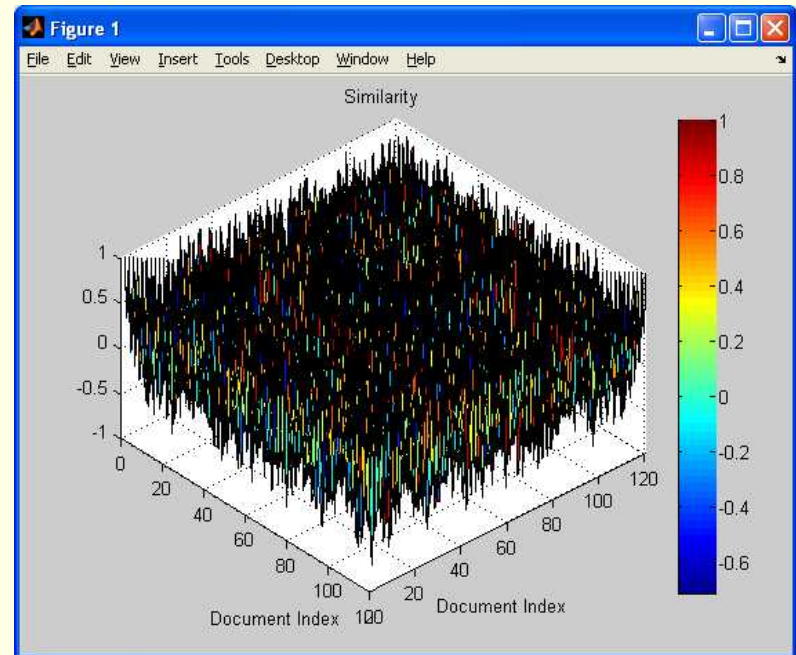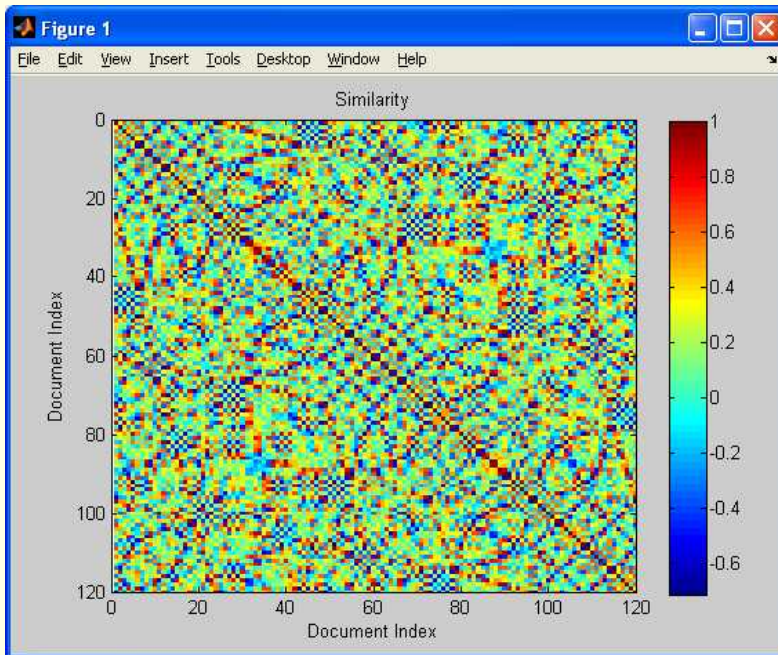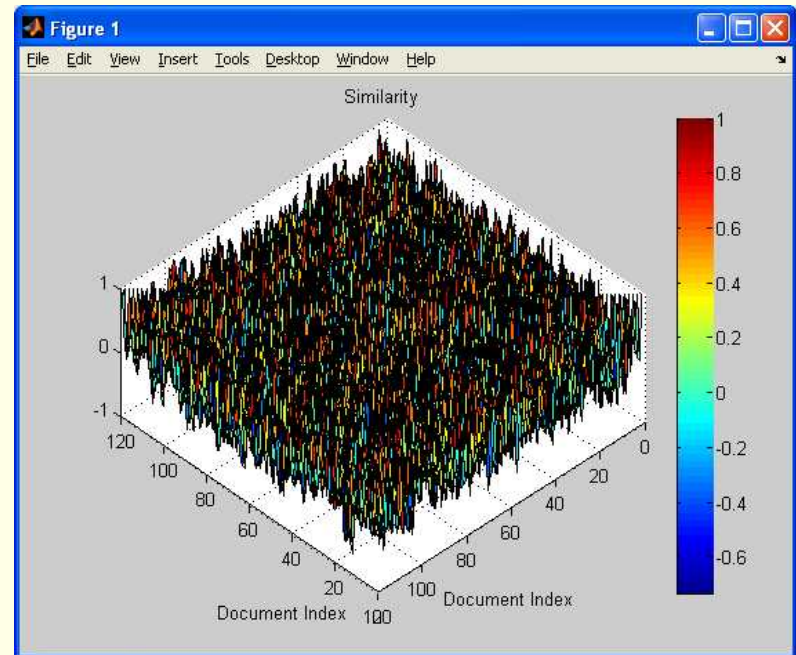
# Document-to-Document Similarities

## • All permutations (n = 2 word-grams) - LSA

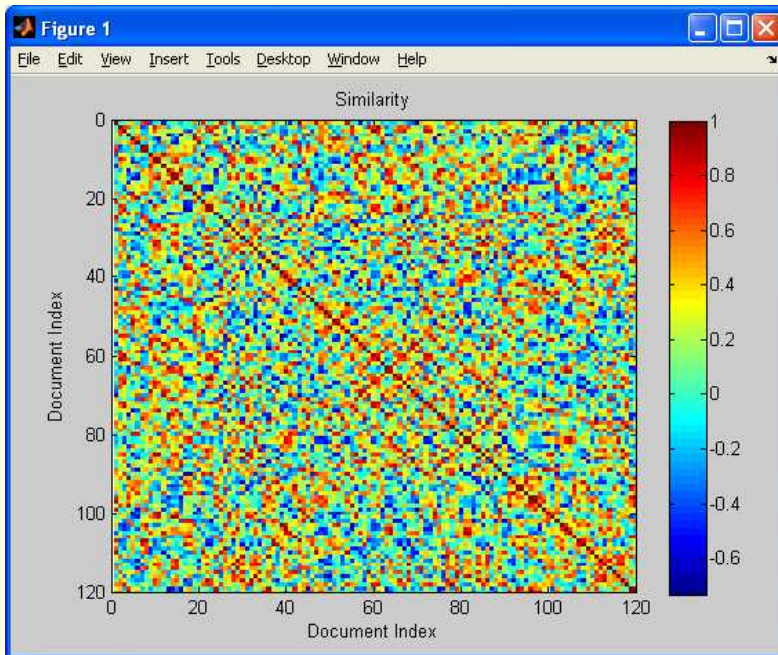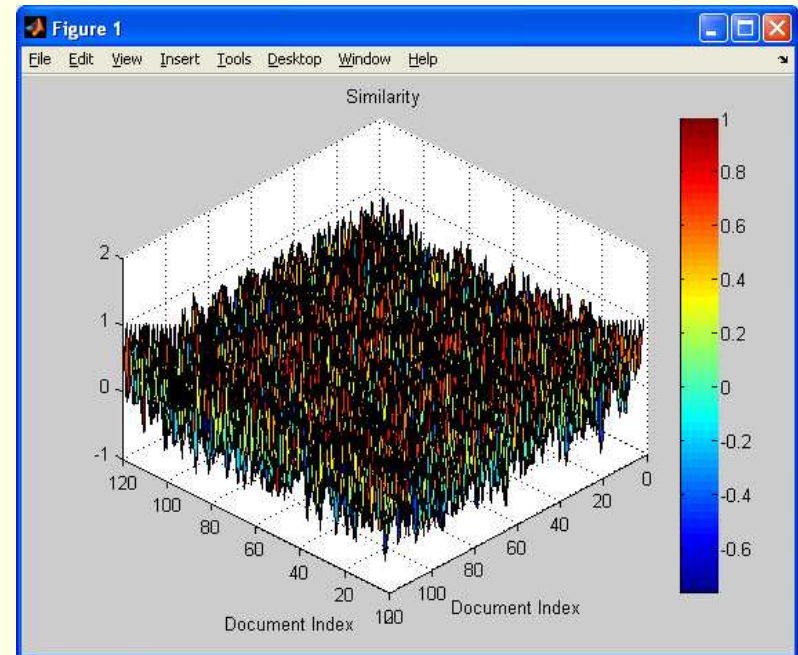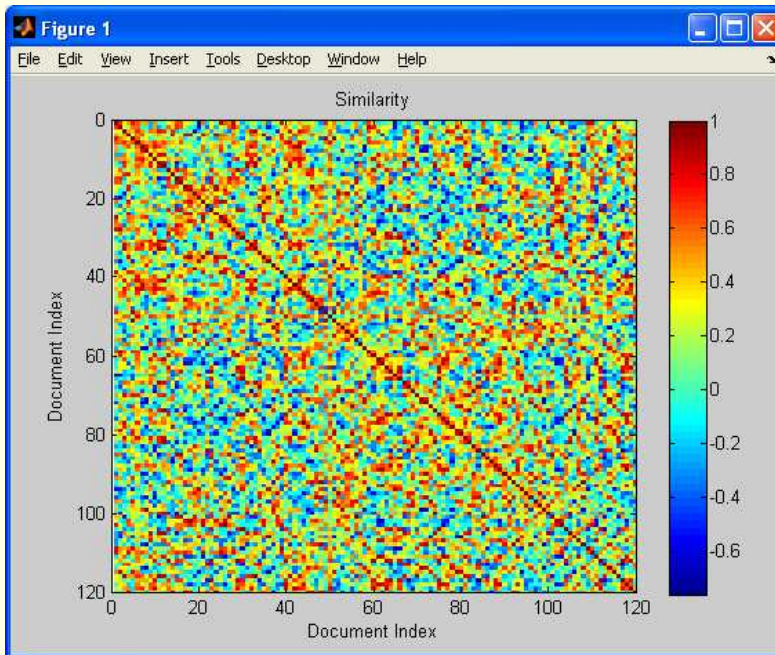# Document-to-Document Similarities

• **All permutations (2 <= n <= 3 word-grams) - LSA**

**Cognitive Systems** Applying computer models of human cognition to create unique technology solutions

# Document-to-Document Similarities

• **All permutations (n <= 5 word-grams) - LSA**





Sandia National Laboratories

# Document-to-Document Similarities

• **All permutations (n = 5 word-grams) - LSA**