# Multilingual Text Analysis of Large Data

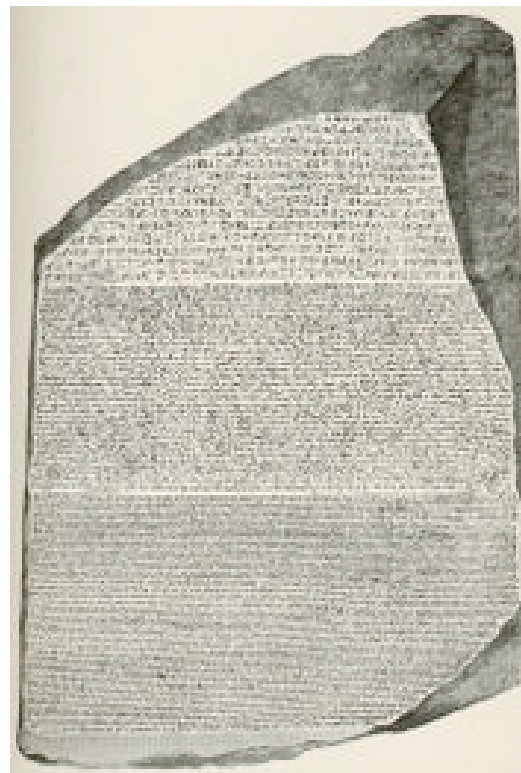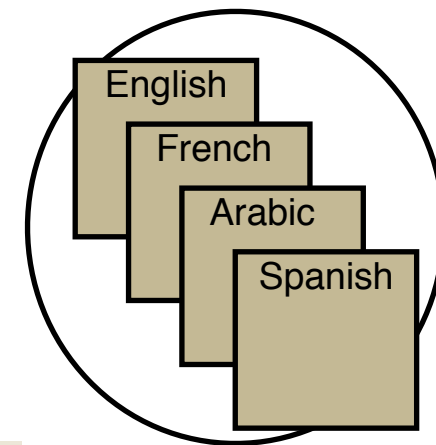**Brett Bader** *and Ron Oldfield*

Sandia National Laboratories

January 20, 2011

# SNL has developed multilingual text analysis to link threats across multiple languages

- "Translate" new documents into a language-independent concept space, which is useful for:
  - Translation triage (i.e., translate documents in clusters of interest)
  - Ideological classification (e.g., hostile to democracy)
  - Multilingual sentiment analysis

Sandia's database: 54 languages: >99% coverage of web

| | | |
|---|---|---|
| Afrikaans | Estonian | Norwegian |
| Albanian | Finnish | Persian (Farsi) |
| Amharic | French | Polish |
| Arabic | German | Portuguese |
| Aramaic | Greek (New Testament) | Romani |
| Armenian Eastern | Greek (Modern) | Romanian |
| Armenian Western | Hebrew (Old Testament) | Russian |
| Basque | Hebrew (Modern) | Scots Gaelic |
| Breton | Hungarian | Spanish |
| Chamorro | Indonesian | Swahili |
| Chinese (Simplified) | Italian | Swedish |
| Chinese (Traditional) | Japanese | Tagalog |
| Croatian | Korean | Thai |
| Czech | Latin | Turkish |
| Danish | Latvian | Ukrainian |
| Dutch | Lithuanian | Vietnamese |
| English | Manx Gaelic | Wolof |
| Esperanto | Maori | Xhosa |

English
French
Arabic
Spanish

Sandia National Laboratories

# Bag of Words/Vector Space Model

## Documents

D1:   How to Bake Bread Without Recipes
D2:   The Classic Art of Viennese Pastry
D3:   Numerical Recipes: The Art of Scientific Computing
D4:   Breads, Pastries, Pies and Cakes: Quantity Baking Recipes
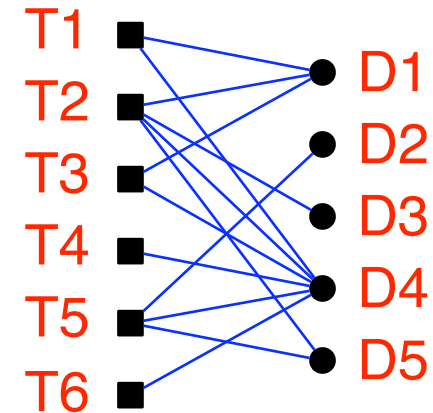D5:   Pastry: A Book of Best French Recipes

## Terms

T1:   bak(e,ing)
T2:   recipes
T3:   bread
T4:   cake
T5:   pastr(y,ies)
T6:   pie

### Key concepts
- Bag of words
- Stemming
- Vector space model
- Scaling for information content
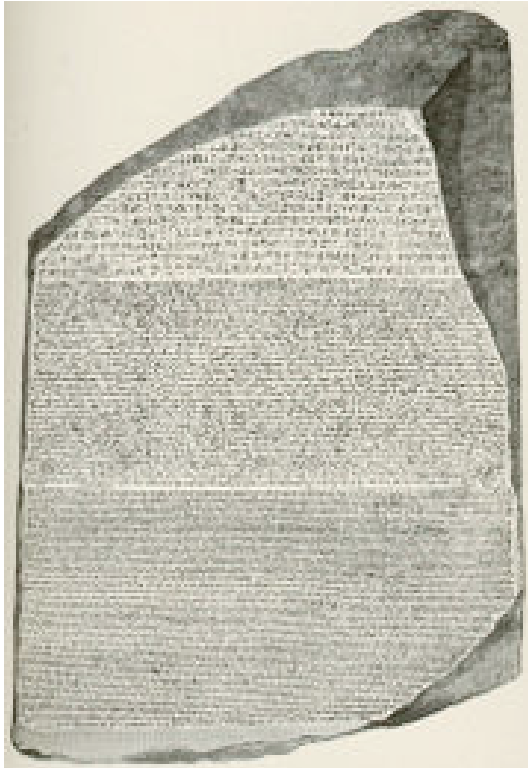
## Bipartite graph

T1
T2
T3
T4
T5
T6

D1
D2
D3
D4
D5

## Term-by-doc (adjacency) matrix

$$\hat{A} = \begin{pmatrix} 1 & 0 & 0 & 1 & 0 \\ 1 & 0 & 1 & 1 & 1 \\ 1 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 & 1 \\ 0 & 0 & 0 & 1 & 0 \end{pmatrix} \begin{matrix} T1 \\ T2 \\ T3 \\ T4 \\ T5 \\ T6 \end{matrix}$$

D1  D2  D3  D4  D5
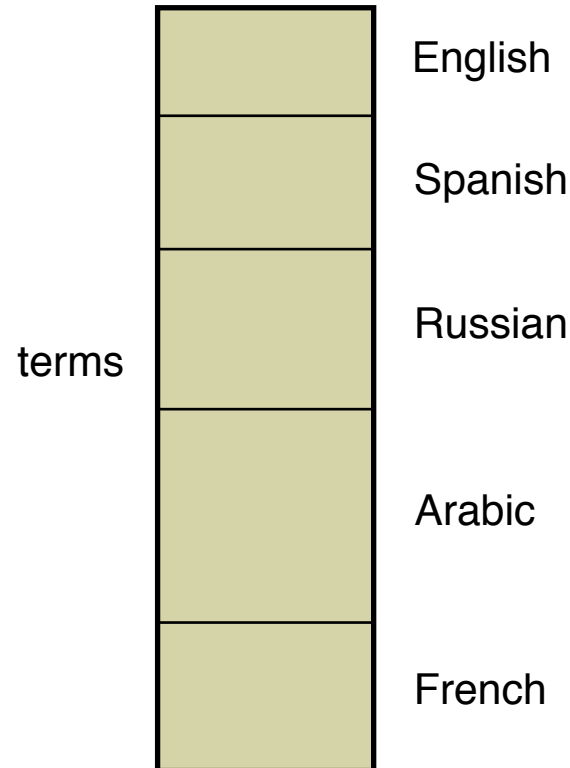
Sandia National Laboratories

# Term-Document Matrix

Term-by-verse matrix
for all languages

Rosetta Stone



Bible verses

terms

English

Spanish

Russian

Arabic

French

163,745 x 31,230

Look for co-occurrence of terms in the same verses and across languages to capture latent concepts

- *Multi-parallel* corpora
  - Bible
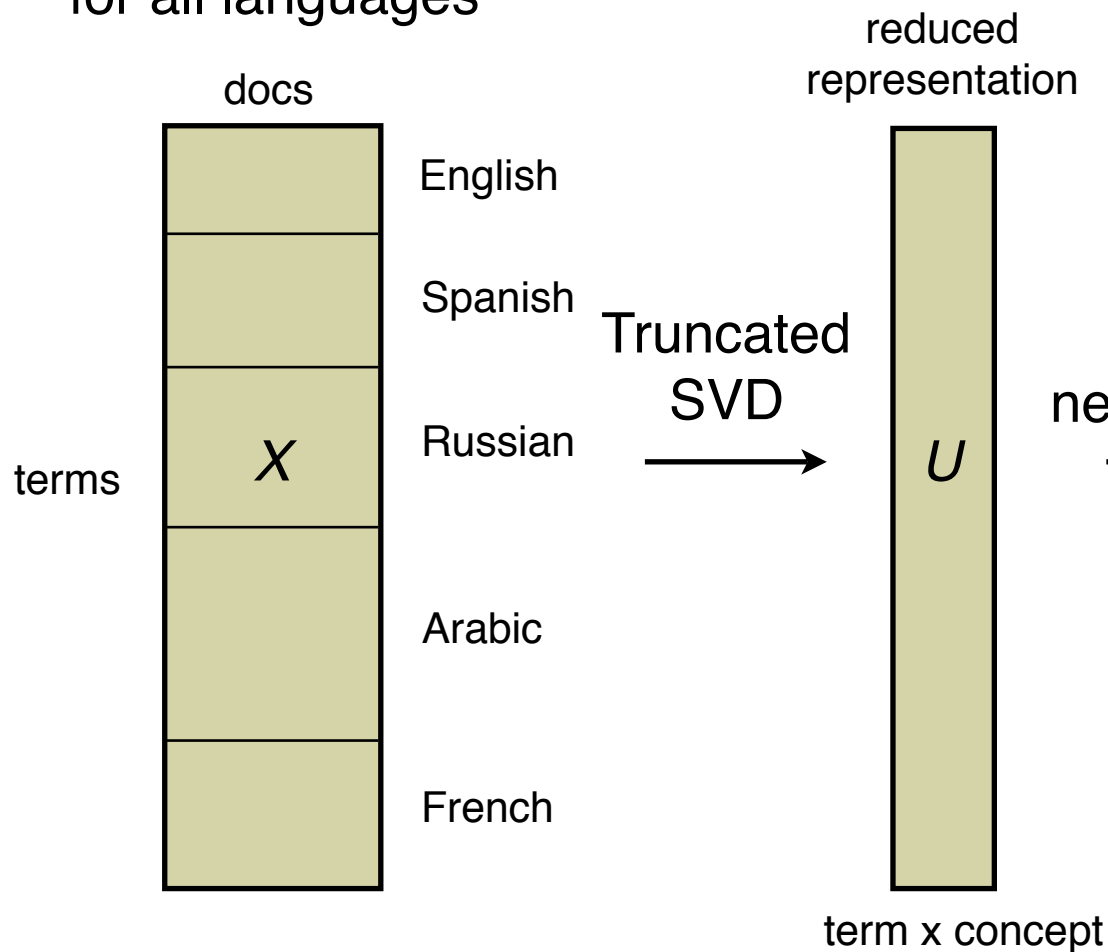  - Quran
  - European Parliament proceedings (Europarl)*

Sandia
National
Laboratories

# Europarl Corpus

- Extracted from the proceedings of the European Parliament
- Translations in 11 languages
  - French, Italian, Spanish, Portuguese (Romantic)
  - English, Dutch, German, Danish, Swedish (Germanic)
  - Greek
  - Finnish
- Sentence aligned text (16 M sentences across 11 languages)
- 1,247,832 speeches (including translations)
- 1,249,253 terms (from all 11 languages)
  - English terms: 46,074

# Multilingual Latent Semantic Analysis
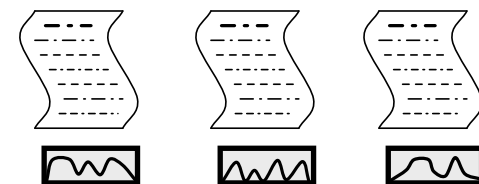
Term-by-doc matrix
for all languages

"Translate" new documents
into a small number of
language-independent features

docs

reduced
representation

| English |
| Spanish |
| Russian |
| Arabic |
| French |

terms

$X$

Truncated
SVD

$U$

term x concept

Project
new documents

| | |
|---|---|
| dimension 1 | 0.1375 |
| dimension 2 | 0.1052 |
| dimension 3 | 0.0341 |
| dimension 4 | 0.0441 |
| dimension 5 | -0.0087 |
| dimension 6 | 0.0410 |
| dimension 7 | 0.1011 |
| dimension 8 | 0.0020 |
| dimension 9 | 0.0518 |
| dimension 10 | 0.0822 |
| dimension 11 | -0.0101 |
| dimension 12 | -0.1154 |
| dimension 13 | -0.0990 |
| dimension 14 | 0.0228 |
| dimension 15 | -0.0520 |
| dimension 16 | 0.1096 |
| dimension 17 | 0.0294 |
| dimension 18 | 0.0495 |
| dimension 19 | 0.0553 |
| dimension 20 | 0.1598 |

Document feature
vector

Applications
- cross-language retrieval
- pairwise similarities for clustering
- machine learning applications

Sandia National Laboratories

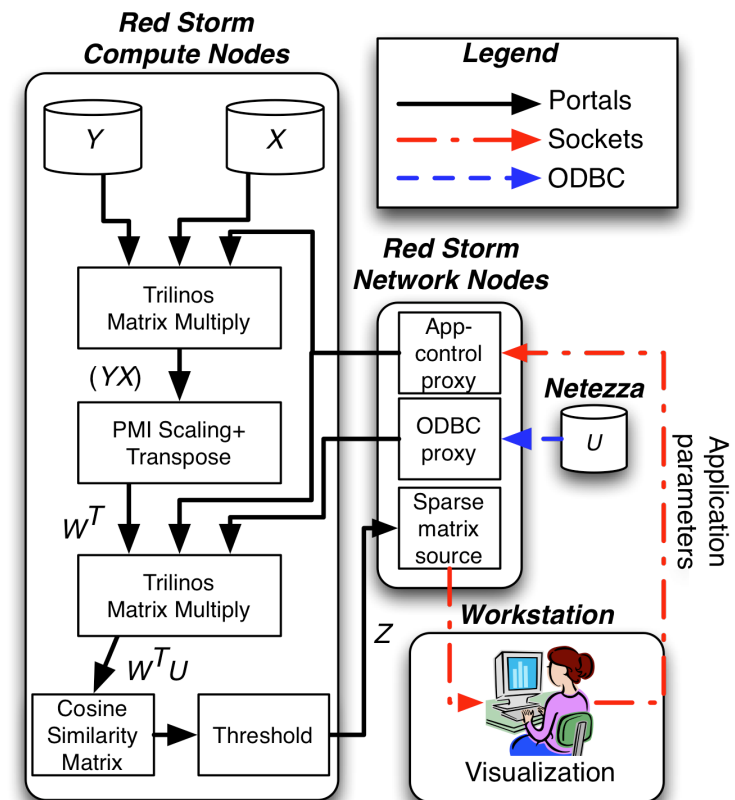# Multilingual Clustering is a
# Great Candidate for HPC

– Scale of Data
  - Millions of elements (Bible, Quran, Wikipedia, Europarl)
  - Computationally expensive (matrix multiplies for large matrices)

– Time to Solution
  - Interactive control/vis is a motivating factor
  - Focus on "strong scaling" capabilities of HPC platform

– Leveraging Existing Sandia Libraries
  - LMSA for dataset generation
  - Trilinos for computation
  - Titan for visualization
  - Nessie for data services (provides "glue" to integrate systems)

Sandia
National
Laboratories

# Architectural Challenges

Exploiting specialized architectures
- Red Storm for numerics
- Clusters/Workstations for vis and interactive control
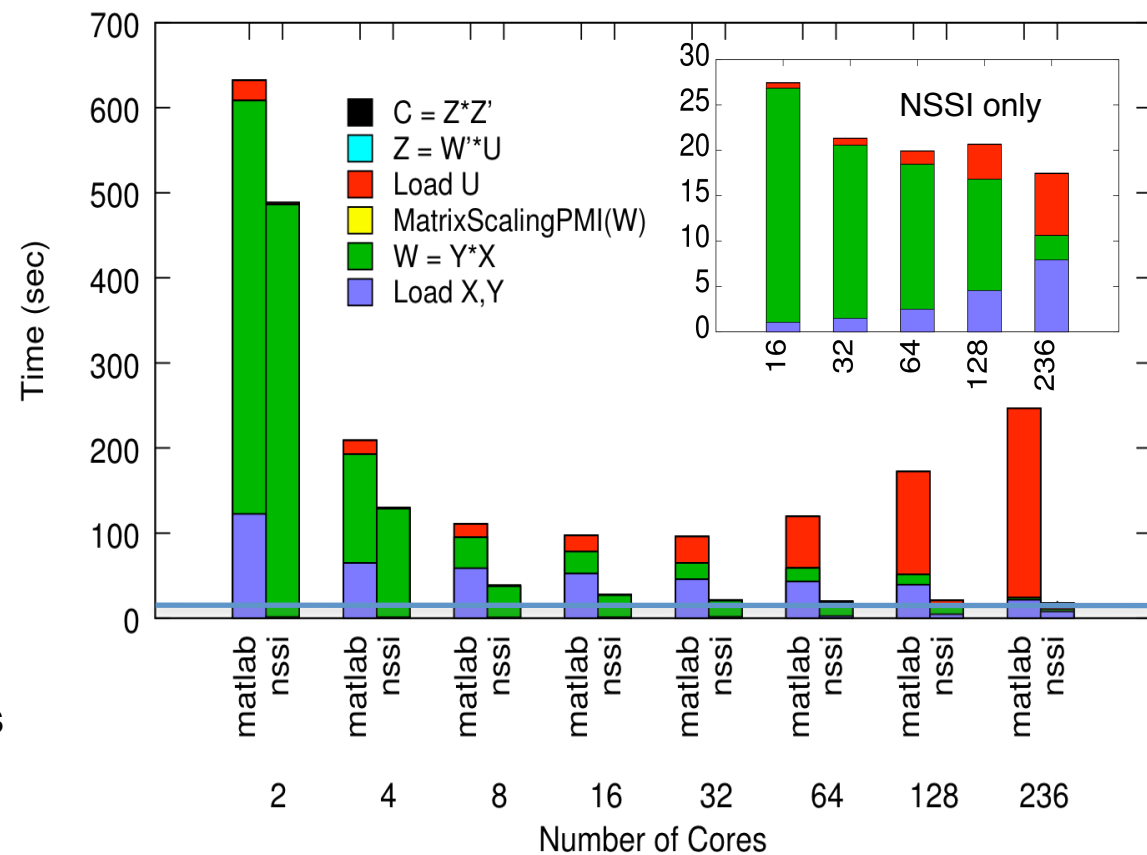- Data Warehouse Appliances for database functionality



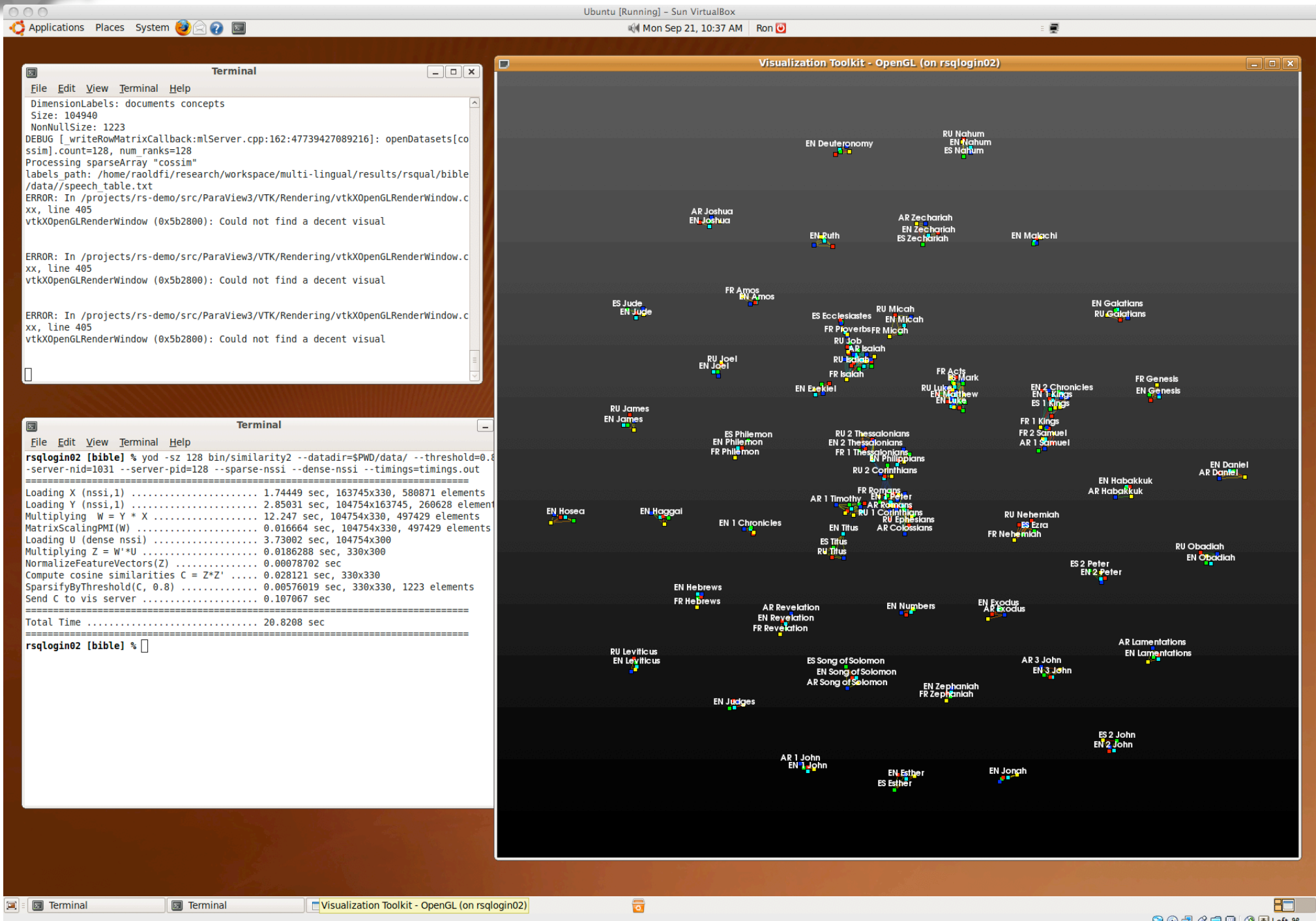*Integrating these systems for interactive jobs has never been done*

# Scaling Challenges for Multilingual Clustering

- **Strong scaling exposes weaknesses in loading**
  - Original methods for loading were not designed for production use.
- **Improvements**
  - Sparse Reads
    - Keep track of processor mapping information
    - Parallel I/O
  - Dense Reads
    - Convert to binary format
    - Parallel I/O
    - Data ordering
- **Status on Red Storm (Cray XT4)**
  - 250K docs of Europarl dataset requires 2048 nodes to execute (memory constrained)
  - At 4096 cores, we overwhelm network communication layer when reading input
  - Our target data set has over 1M docs



Performance Results: Bible Dataset

Legend:
- C = Z*Z'
- Z = W'*U
- Load U
- MatrixScalingPMI(W)
- W = Y*X
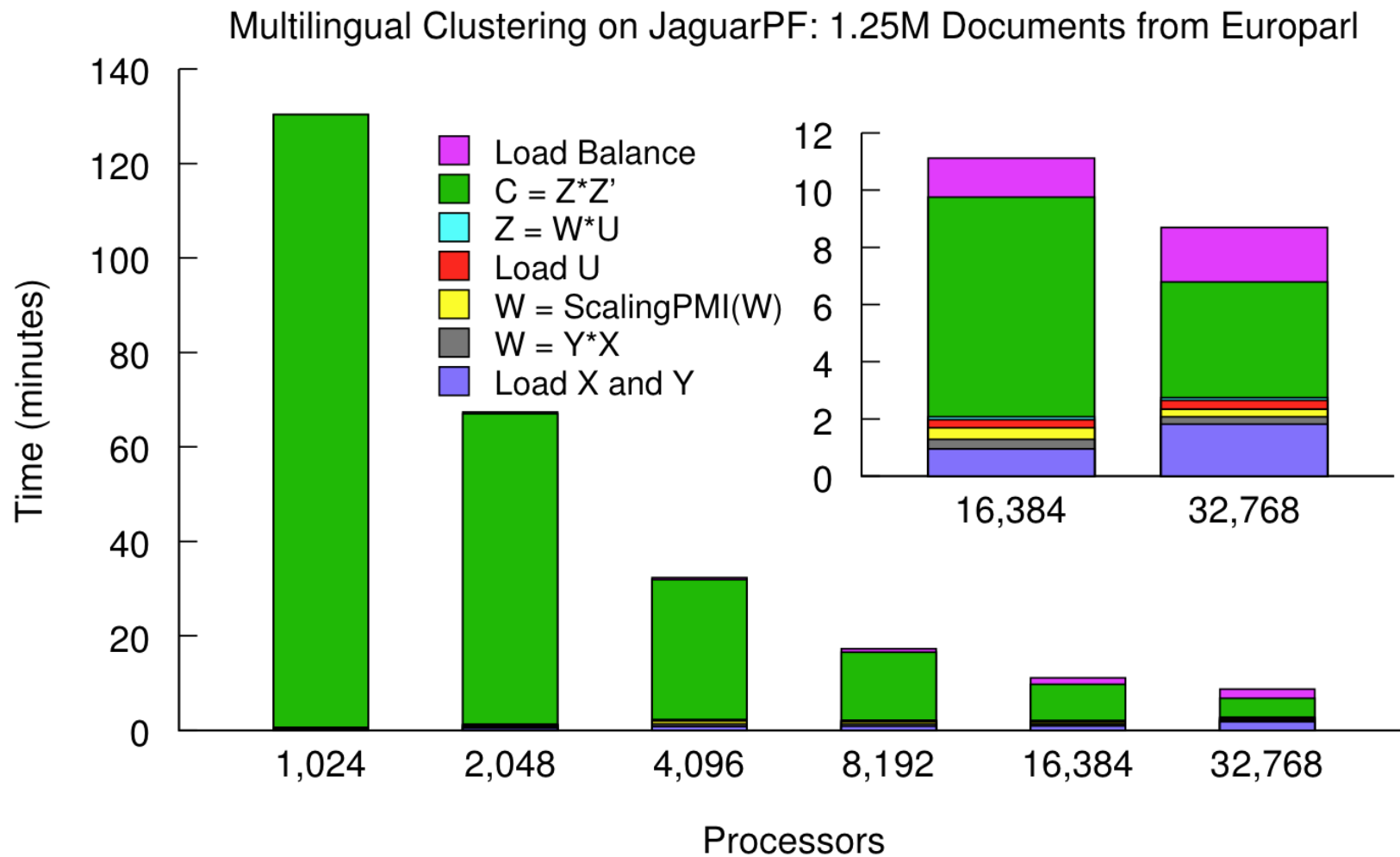- Load X,Y

NSSI only

Sandia National Laboratories

# HPC clustering

# Scaling Challenges for Multilingual Clustering

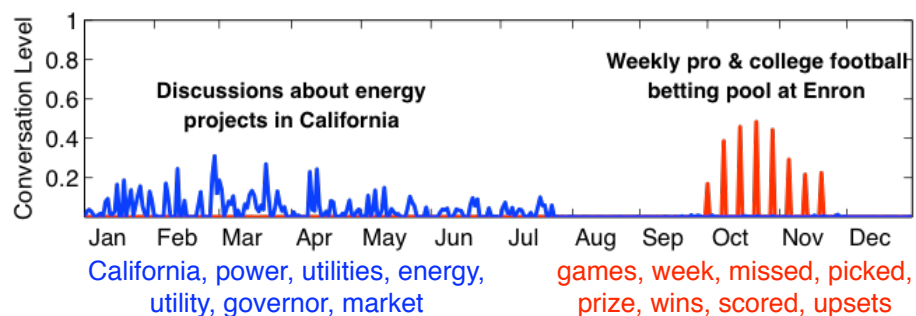- Performance on JaguarPF (Cray XT5)
  - 1.25M docs of Europarl data set
  - With 32K cores, it takes 470 seconds

Multilingual Clustering on JaguarPF: 1.25M Documents from Europarl

Legend:
- Load Balance
- C = Z*Z'
- Z = W*U
- Load U
- W = ScalingPMI(W)
- W = Y*X
- Load X and Y

Sandia National Laboratories
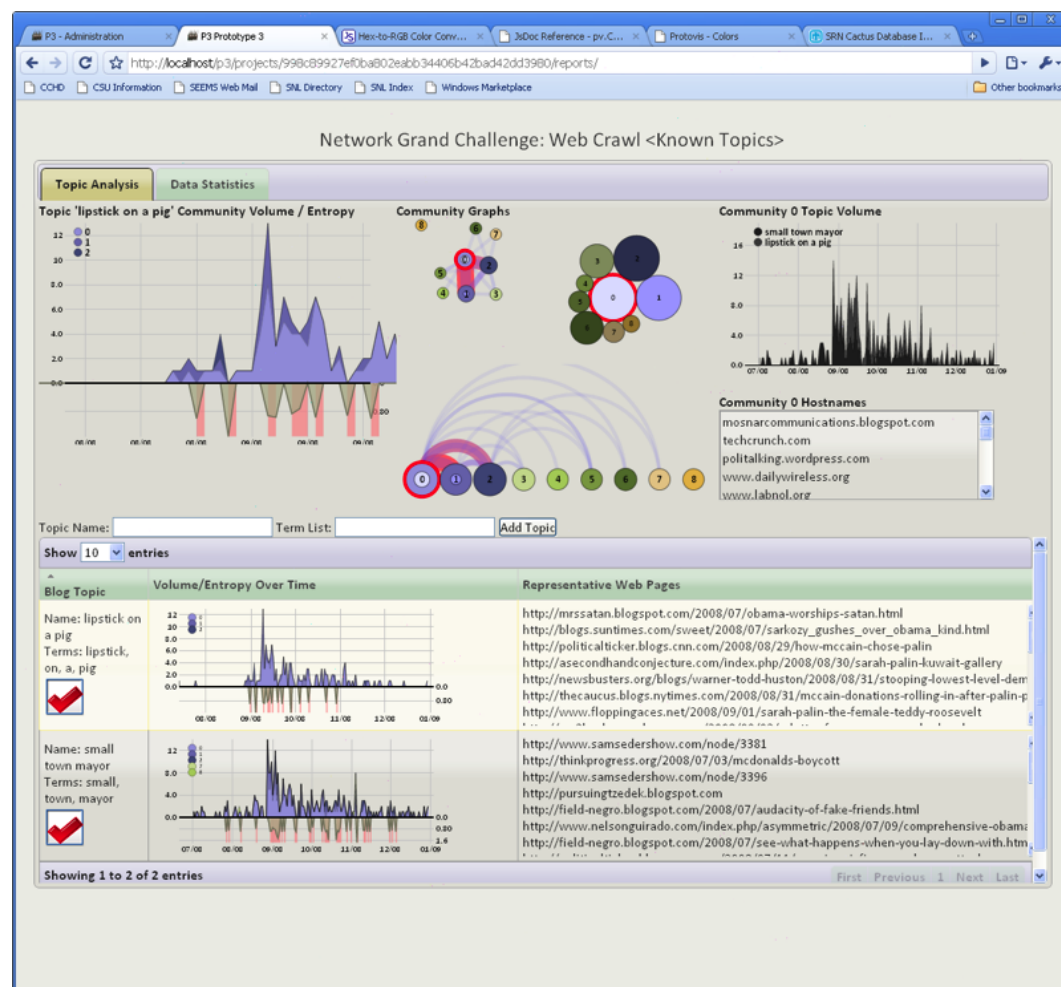
# Related Text Analysis Projects

- Discussion tracking in Enron emails
- Uncovering plots in text (scenario discovery)
- Network data exfiltration analysis
- Higher-order web link analysis
- Unsupervised part-of-speech tagging
- Multilingual sentiment analysis
- Identifying emerging keywords of interest

*Analysis tools for web forecasting*

*Identifying unusual activity in Enron emails*



California, power, utilities, energy, utility, governor, market

games, week, missed, picked, prize, wins, scored, upsets

J. Steffes, S. Kean, J. Dasovich, R. Shapiro, P. Allen, ...

A. Pace, L. Campbell, C. Dean

# Selected References

- US Patent Application SD11033, "Technique for Information Retrieval Using Enhanced Latent Semantic Analysis," Peter Chew and Brett Bader, filed 2009.

- Bader, Kegelmeyer, and Chew. 2011. Multilingual sentiment analysis using latent semantic indexing and machine learning. Submitted to ACL.

- Bader and Chew. 2008. Enhancing multilingual latent semantic analysis with term alignment information. *COLING 2008.*

- Chew, Bader, and Abdelali. 2008. Latent Morpho-Semantic Analysis: Multilingual Retrieval with character N-grams and mutual information. *COLING 2008*.

- Chew, Kegelmeyer, Bader and Abdelali. 2008. The Knowledge of Good and Evil: Multilingual Ideology Classification with PARAFAC2 and Machine Learning. *Language Forum* (34), 37-52.

- Chew, Bader, Kolda and Abdelali. 2007. Cross-language information retrieval using PARAFAC2. *Proceedings of KDD 2007.*

- Chew and Abdelali. 2007. Benefits of the 'massively parallel Rosetta Stone': cross-language information retrieval with over 30 languages, *Proceedings of the Association for Computational Linguistics conference,* 2007.

- Chew, Verzi, Bauer and McClain. 2006. Evaluation of the Bible as a Resource for Cross-Language Information Retrieval. *Proceedings of the Workshop on Multilingual Language Resources and Interoperability,* 2006, 68–74.

Brett Bader (bwbader@sandia.gov)
http://www.sandia.gov/~bwbader

Sandia National Laboratories