# Bilingual Document Alignment with Latent Semantic Indexing
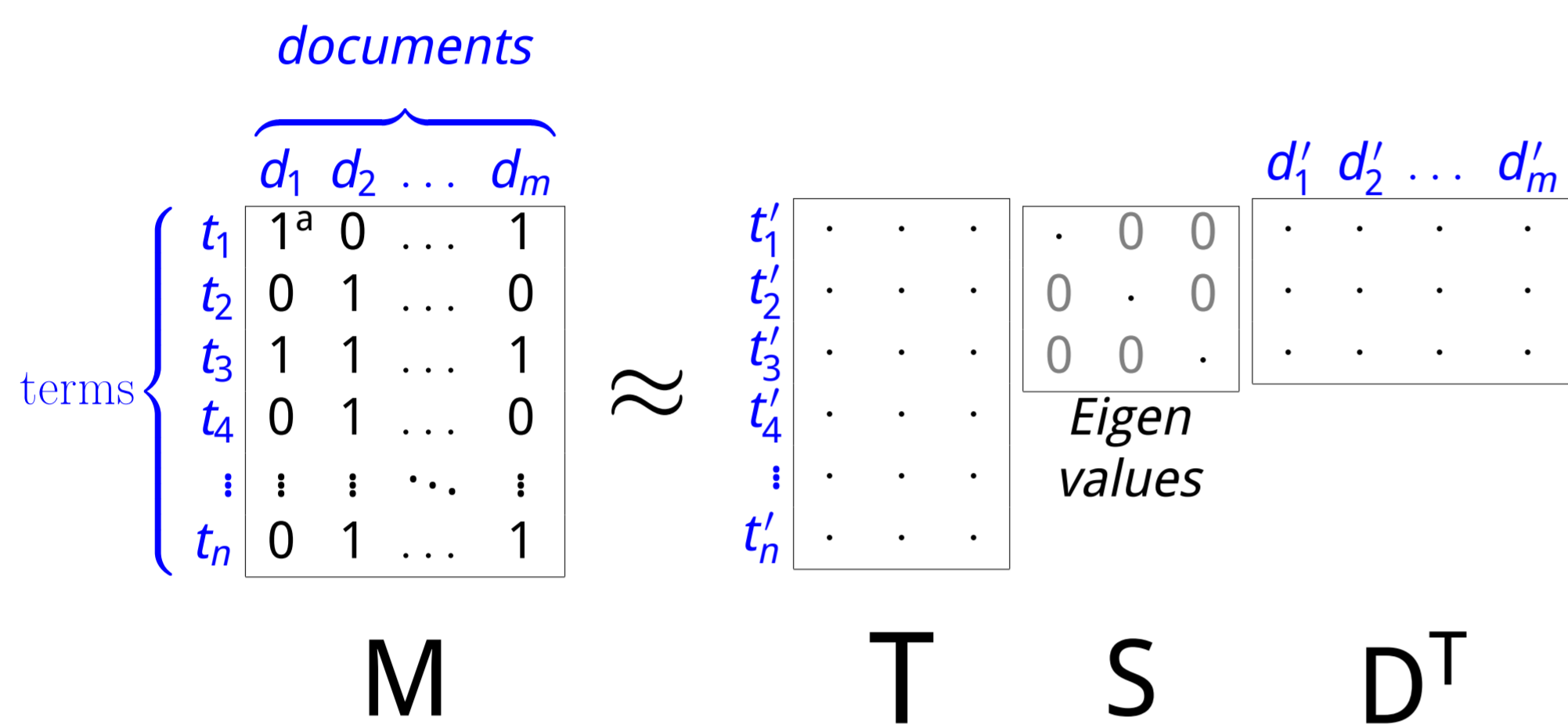## Ulrich Germann, University of Edinburgh

## The Bilingual Document Alignment Task

Given a collection of web page downloads and their corresponding ULRs, find pairs of pages that are translations of one another. A seed set of page pairs is given for development purposes.

## Basic Approach

Compute the cosine between **document embeddings** in a joint semantic space to measure cross-lingual document similarity. We use *Latent Semantic Indexing* (LSI) to map documents into the joint space. LSI relies on reduced-rank *Singular Value Decomposition* (SVD) to perform this task.



[a] in practice, we use tf-idf weighted term counts; see below.

## Technical Procedure

1. Fill the term-document matrix with $tf_{t,d} \cdot idf_{t,D_d}$ values for each term occurring in each document pair $d$ from the training data, with

$$tf_{t,d} = 1 + log(f_{t,d})$$

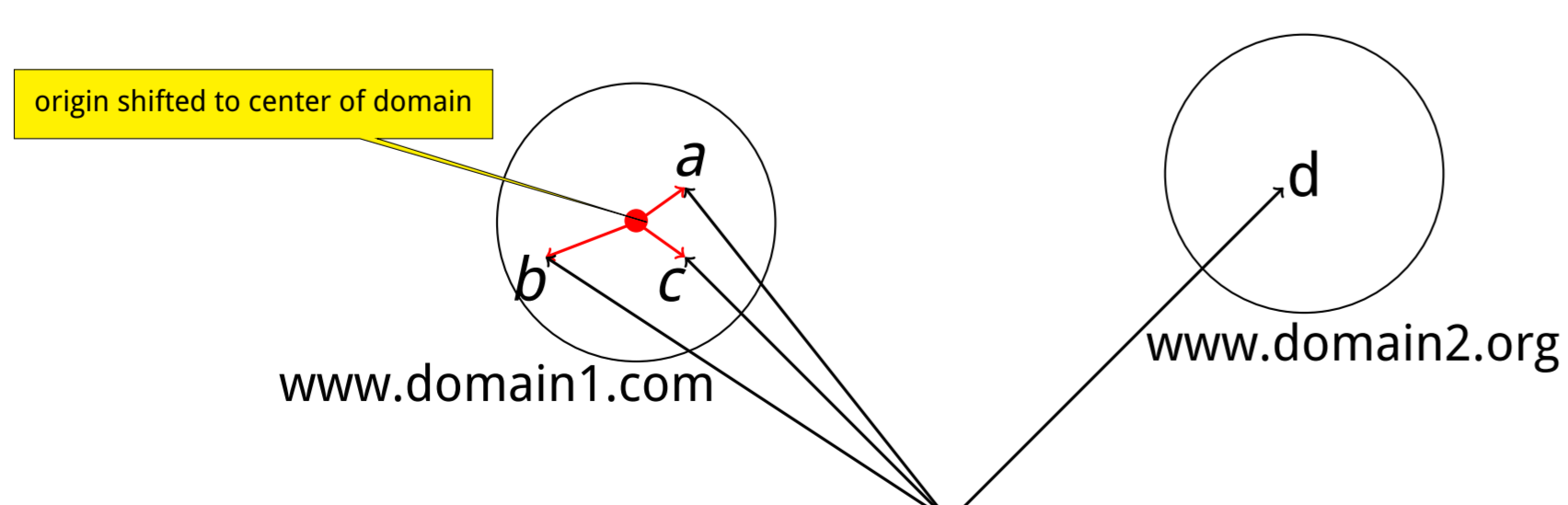$$idf_{t,D_d} = \log \frac{|D_d|}{|\{d' \in D : t \in d'\}|}$$

where $f_{t,d}$ is the raw occurrence count of term $t$ in $d$ and $D_d$ the web domain that document $d$ belongs to.

2. Factorize term-document matrix with SVD.
3. For each document in the data set, compute a document vector $\vec{d}$ as in Step 1.
4. Fold it in: $\vec{d}' = S^{-1} T^{\mathsf{T}} \vec{d}$.
5. Compute document similarity between two documents $d_x$ an $d_y$ as

$$sim(d_x, d_y) = \cos(S \vec{d}'_x, S \vec{d}'_y).$$

## Features

- URL string similarity (details in the paper)
- Cosine in the global space (cos)
- Cosine after shifting origin to center of documents from a specific web site (lcos) by subtracting the mean document vector over all documents from that site from each individual vector



## Results (rank=1000)

### recall on the training data

| features used | | | seed data included | | | | seed data excluded | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | strict[a] | 1.00[b] | 0.99[b] | 0.95[b] | 0.90[b] | strict[a] | 1.00[b] | 0.99[b] | 0.95[b] | 0.90[b] |
| cosine (cos) | | | 86.7 | 93.4 | 95.4 | 96.7 | 97.6 | 82.5 | 88.9 | 91.3 | 92.9 | 93.7 |
| "local" cos. (lcos) | | | 86.7 | 92.8 | 94.7 | 95.8 | 96.9 | 83.3 | 88.9 | 91.4 | 92.8 | 93.6 |
| URL similarity (url) | | | 83.6 | 87.8 | 88.1 | 88.2 | 88.2 | 83.6 | 87.8 | 88.1 | 88.2 | 88.2 |
| cos | lcos | | 87.2 | 93.7 | 95.6 | 96.6 | 97.5 | 83.3 | 89.7 | 92.1 | 93.6 | 94.4 |
| cos | | url | 90.6 | 94.7 | 95.6 | 96.4 | 97.1 | 86.3 | 90.6 | 91.4 | 92.7 | 93.5 |
| | lcos | url | 91.3 | 95.4 | 96.3 | 97.2 | 97.8 | 86.8 | 91.3 | 92.2 | 93.4 | 94.2 |
| cos | lcos | url | 92.8 | 96.7 | 97.6 | 98.5 | 99.1 | 88.0 | 92.5 | 93.4 | 94.7 | 95.5 |

### recall on the test data

| features used | | | strict[a] | 1.00[b] | 0.99[b] | 0.95[b] | 0.90[b] |
|---|---|---|---|---|---|---|---|
| cos | lcos | url | | | 87.6 | 87.6 | 94.1 | 95.5 | 96.0 |

[a] exact string match with the reference ULR pairs
[b] soft match based on document similarity with different thresholds.

## Special Characteristics of Web Pages

- Web pages are a mix of **boilerplate text** and **payload**. **Boilerplate text** (menus, disclaimers, etc.) occurs on most web pages that belong to a particular (sub-)domain. The **payload** is the actual "content" of the page. Boilerplate text is very specific to specific domains but provides little information about the payload.
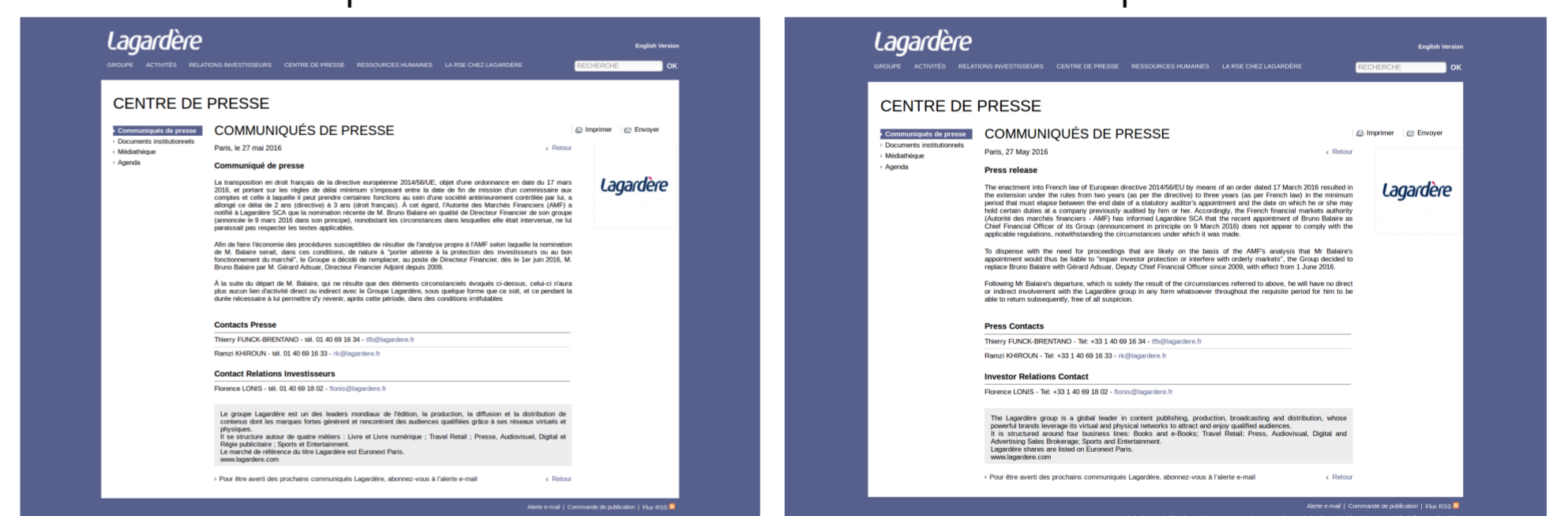


- Different URLs may lead to the same (or almost the same) page. This leads to problems if they rely only on pairs of URLs for evaluation.
- Web pages are often dynamically generated, sometimes delivering the same payload with different "facades" (e.g., web view, print view, boilerplate in different languages).