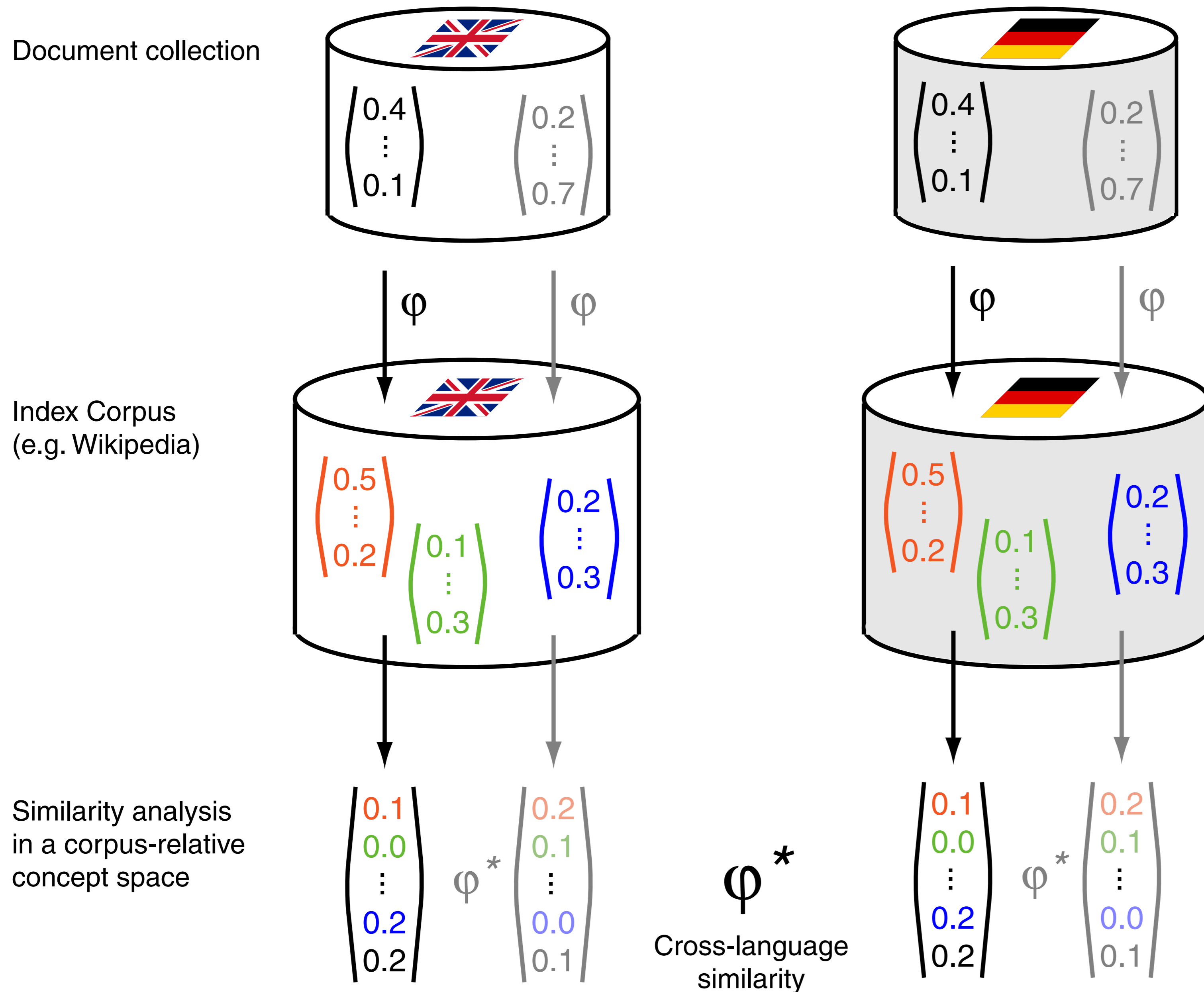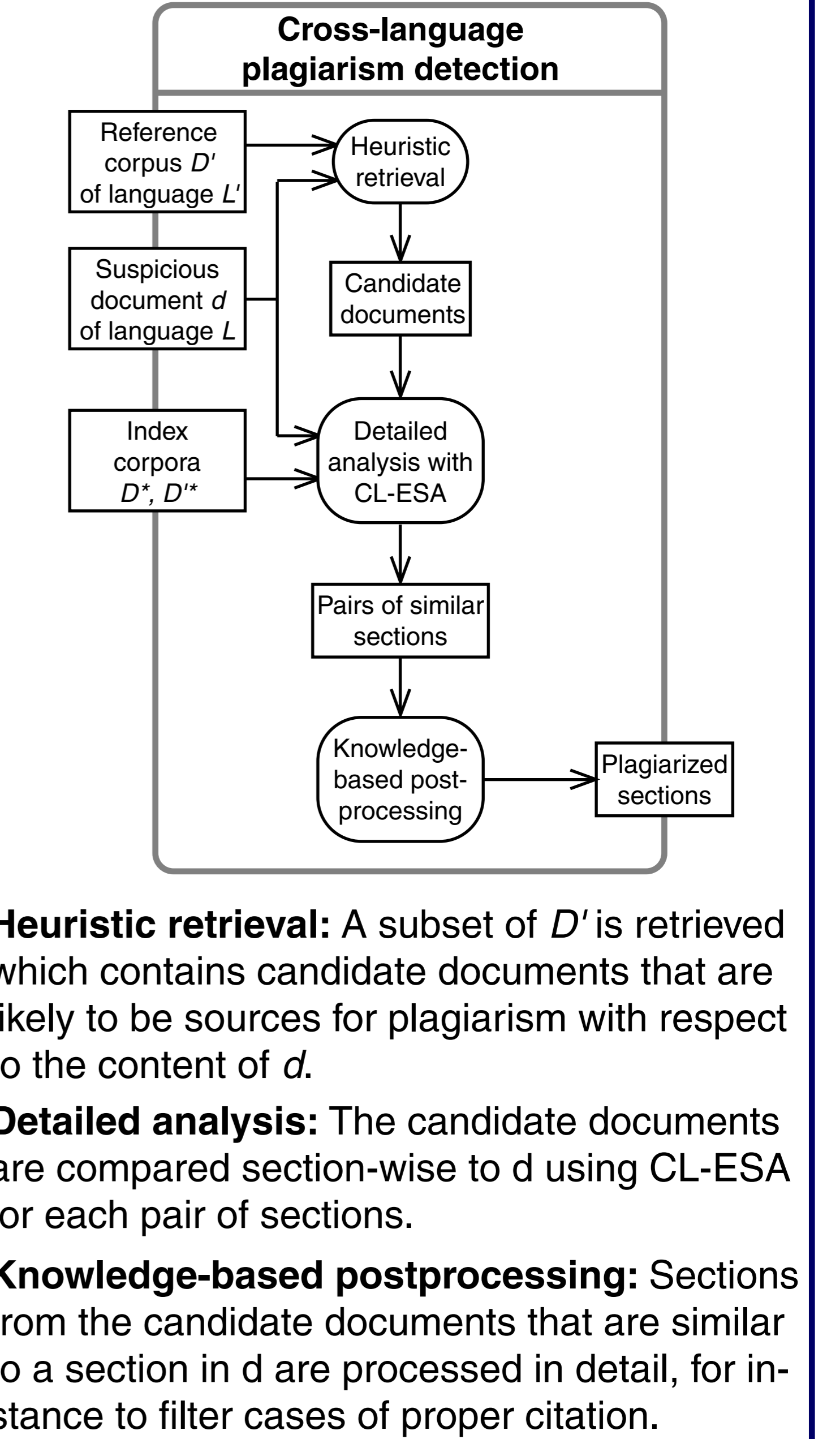# A Wikipedia-based Multilingual Retrieval Model

Martin Potthast, Benno Stein, and Maik Anderka
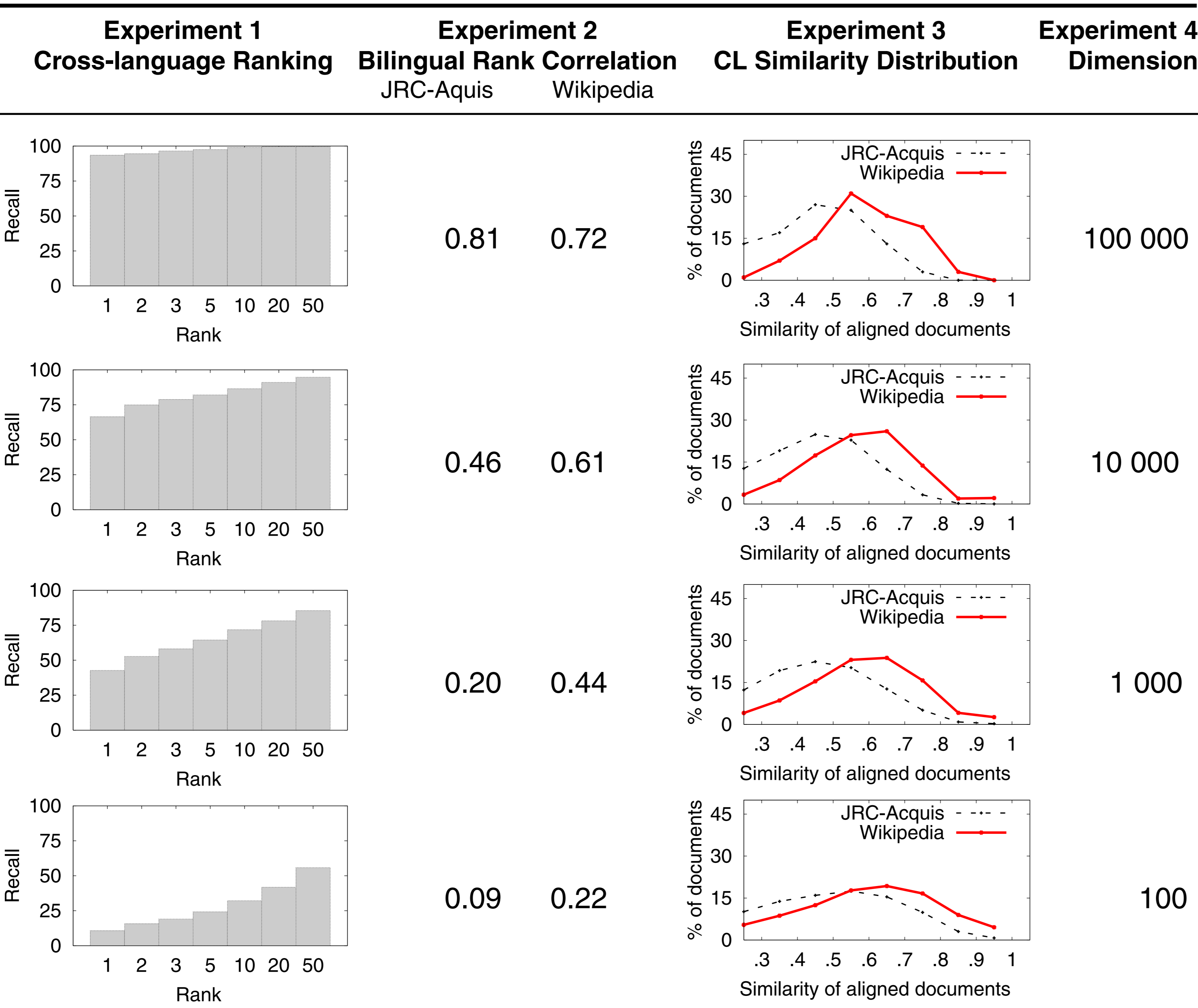
## Contribution

Document collection

0.4 ⋮ 0.1    0.2 ⋮ 0.7

0.4 ⋮ 0.1    0.2 ⋮ 0.7

φ    φ    φ    φ

Index Corpus (e.g. Wikipedia)

0.5 ⋮ 0.2    0.1 ⋮ 0.3    0.2 ⋮ 0.3

0.5 ⋮ 0.2    0.1 ⋮ 0.3    0.2 ⋮ 0.3

Similarity analysis in a corpus-relative concept space

0.1 0.0 ⋮ 0.2 0.2    φ*    0.2 0.1 ⋮ 0.0 0.1

φ*    Cross-language similarity

0.1 0.0 ⋮ 0.2 0.2    φ*    0.2 0.1 ⋮ 0.0 0.1

## Current Work

**Cross-language plagiarism detection**

Reference corpus $D'$ of language $L'$ → Heuristic retrieval

Suspicious document $d$ of language $L$ → Candidate documents

Index corpora $D^*$, $D'^*$ → Detailed analysis with CL-ESA

Pairs of similar sections

Knowledge-based post-processing → Plagiarized sections

**Heuristic retrieval:** A subset of $D'$ is retrieved which contains candidate documents that are likely to be sources for plagiarism with respect to the content of $d$.

**Detailed analysis:** The candidate documents are compared section-wise to d using CL-ESA for each pair of sections.

**Knowledge-based postprocessing:** Sections from the candidate documents that are similar to a section in d are processed in detail, for instance to filter cases of proper citation.

## Evaluation

|  | Experiment 1 Cross-language Ranking | Experiment 2 Bilingual Rank Correlation | | Experiment 3 CL Similarity Distribution | Experiment 4 Dimension |
|---|---|---|---|---|---|
|  |  | JRC-Aquis | Wikipedia |  |  |
|  | *(recall vs rank bar chart)* | 0.81 | 0.72 | *(similarity distribution plot)* | 100 000 |
|  | *(recall vs rank bar chart)* | 0.46 | 0.61 | *(similarity distribution plot)* | 10 000 |
|  | *(recall vs rank bar chart)* | 0.20 | 0.44 | *(similarity distribution plot)* | 1 000 |
|  | *(recall vs rank bar chart)* | 0.09 | 0.22 | *(similarity distribution plot)* | 100 |

Experiment 1 charts: Recall (100, 75, 50, 25, 0) vs Rank (1, 2, 3, 5, 10, 20, 50)

Experiment 3 charts: JRC-Acquis (dashed) and Wikipedia (red); % of documents (45, 30, 15, 0) vs Similarity of aligned documents (.3 .4 .5 .6 .7 .8 .9 1)

**Test documents:** Two test document collections, $D$ and $D'$, comprising 3 000 documents each were chosen (1 000 translation-aligned, 1 000 concept-aligned, 1 000 not aligned).
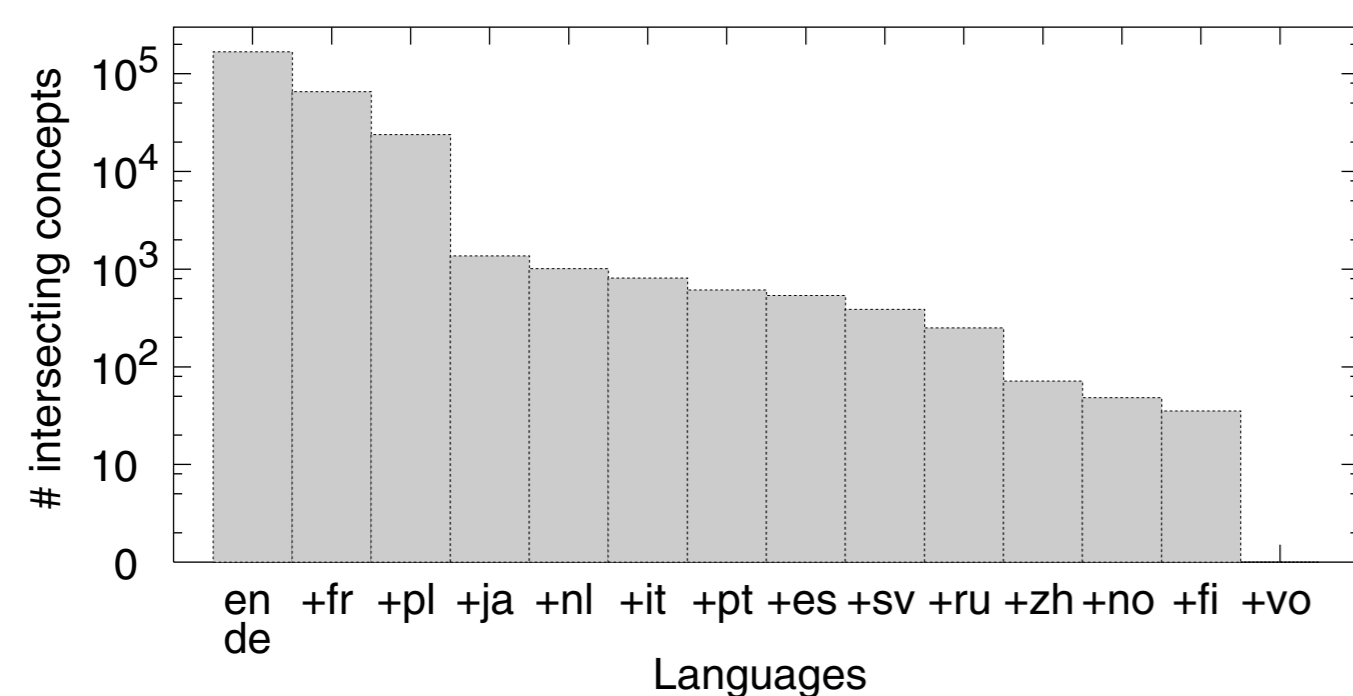
**Experiment 1:** Given an aligned document $d \in D$, all documents in $D'$ are ranked according to their cross-language similarity to $d$. Let $d' \in D'$ be the aligned document of $d \in D$, then the retrieval rank of $d'$ is recorded. Ideally, $d'$ should be on the first or at least on one of the top ranks. The experiment was repeated for all of the aligned documents in $D$.

**Experiment 2:** For a pair of aligned documents $d \in D$ and $d' \in D'$ the documents from D are ranked twice: *(i)* with respect to their cross-language similarity to $d$ using a cross-language retrieval model, and, *(ii)* with respectto their monolingual similarity to $d$ using the vector space model. The top 100 ranks of the two rankings are compared using a rank correlation coefficient, e. g. Spearma's $\rho$, which measures their (dis-)agreement as a value between -1 and 1 respectively.

**Experiment 3:** This experiment contrasts the distribution of pairwise similarities of translation-aligned and concept-aligned documents.

**Experiment 4:** Both retrieval quality and runtime depend on the concept space dimension of CL-ESA, which in turn corresponds the size of a language's index corpus.

*(Experiment 5 chart: # intersecting concepts ($10^5$ to 0) vs Languages: en/de, +fr, +pl, +ja, +nl, +it, +pt, +es, +sv, +ru, +zh, +no, +fi, +vo)*

**Experiment 5:** Starting with the two most prominent languages in Wikipedia, English and German, we study how many concepts are described in both languages, and how many are in the intersection set if more languages are considered.

*(Experiment 6 chart: Time (ms) ($10^4$ to 0) vs Dimensions (10, $10^2$, $10^3$, $10^4$, $10^5$); External memory μ, Internal memory μ, VSM)*

**Experiment 6:** The time to index a document is between 10 to 100 milliseconds, which is comparable to the time to compute a vector space representation. Employed hardware: Intel Core 2 Duo processor at 2 GHz and with 1 GB RAM.