# Candidate Document Retrieval for Web-scale Text Reuse Detection

Matthias Hagen     Benno Stein

Bauhaus-Universität Weimar
matthias.hagen@uni-weimar.de

SPIRE 2011
Pisa, Italy
October 19, 2011

# Text reuse?

Text from one document used in another.

**Mirror**

## Big Ben turning into London's version of the Leaning Tower of Pisa

by Martin Fricker, Daily Mirror 10/10/2011

Big Ben and The Leaning Tower of Pisa (pics: Reuters)

BIG Ben is turning into our own Leaning Tower of Pisa, a worrying survey has confirmed.

The much-loved landmark's tilt has become so pronounced it is noticed by passers-by and tourists.

The Palace of Westminster's clock tower has not been perfectly vertical for years because of shifting ground conditions and tunnelling for Tube lines.

Now engineers say it will one day topple over if the lean is left unchecked. Big Ben is the nickname of the tower's largest bell but the public generally use it as the name of the whole clock, built in 1853.

The peak of the 315ft tower is 18 inches off where it would be if vertical – a 0.26 degree tilt to the north west.

That is one sixteenth of the Pisa tower's lean.

But a survey for London Underground and the Parliamentary Estates Department found the rate of movement accelerated in recent years.

It has caused cracks to appear in walls inside the House of Commons.

Prof John Burland, of Imperial College London, said: "I have heard tourists saying, 'I don't think it is really vertical'. They are quite right. The tilt is now just about visible. If it started greater acceleration we would have to do something in a few years."

The clock moved an eighth of an inch from the perpendicular between November 2002 and August 2003. Since then the tilt has increased 0.04 of a degree each year. At that rate it would crash into Portcullis House, used as MPs' offices – in 5,000 years.

**REUTERS**

## Leaning tower of London? Big Ben is tilting

LONDON | Tue Oct 11, 2011 8:12am EDT

(Reuters) - British landmark Big Ben is leaning to such an extent that the tilt can now be clocked with the naked eye, according to a report commissioned by London Underground and the Parliamentary Estates Department.

The 96 metre (yards) high clock tower of the Houses of Parliament -- known colloquially as Big Ben, the name of the great bell it houses -- is sinking unevenly into the ground, causing it to lean towards the northwest.

"The tilt is now just about visible. You can see it if you stand on Parliament Square and look east, towards the river. I have heard tourists there taking photographs saying 'I don't think it is quite vertical' - and they are quite right," emeritus professor and senior research investigator at Imperial College, London, John Burland, told the Sunday Telegraph.

The level of the tilt has accelerated since 2003, increasing to 0.9 mm a year, compared to the long-term average rate of 0.65 mm a year, the report revealed.

These levels are not considered to be unsafe.

"If it started greater acceleration, we would have to look at doing something but I don't think we need to do anything for a few years yet," Burland said.

Years of underground developments have contributed to the clock tower's tilt, according to the report.

This includes the construction of an underground car park in the early 70s and an extension of the London Underground Jubilee Line, as well as changes in ground conditions.

The tilt has resulted in the formation of cracks in the walls and ceilings of parts of the House of Commons, including the Minister's Wing.

The Palace of Westminster, also known as the Houses of Parliament, is the site of Britain's House of Lords and the House of Commons.

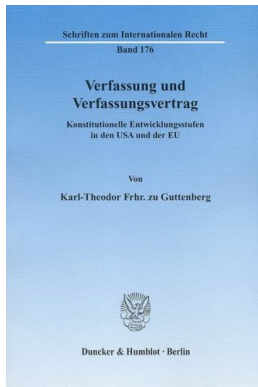The construction of the great clock tower was completed in 1858.

Karl-Theodor zu Guttenberg
(former German Minister of Defence)



60% of dissertation plagiarized

# Paper versions



SPIRE 2011 full paper

# Paper versions

SPIRE 2011 full paper

ECDL 2010 poster

# Text reuse detection

Given        "suspicious" document

# Text reuse detection

Given      "suspicious" document

Step 1:    Find a set of candidate documents

# Text reuse detection

Given      "suspicious" document

Step 1:      Find a set of candidate documents
Step 2:      In-depth analysis against each candidate

We focus on Step 1

Candidate document retrieval

# Candidate document retrieval

## Observations

- Text reuse source      = the entire Web      $\rightarrow$ web search
- Same topic doc's       = more likely source

# Candidate document retrieval

## Observations

- Text reuse source = the entire Web → web search
- Same topic doc's = more likely source

- Too many candidates = bad runtime → system capacity $k$
- Up to $k$ candidates = reasonable runtime

# Candidate document retrieval

## Observations

- Text reuse source    = the entire Web     $\rightarrow$ web search
- Same topic doc's     = more likely source

- Too many candidates = bad runtime     $\rightarrow$ system capacity $k$
- Up to $k$ candidates    = reasonable runtime

## Idea

Retrieve a feasible number of similar web documents.

# Standing on the shoulders of . . .

Random string as query    [Dasdan et al., CIKM 2009]

Rare keywords as query    [Dasdan et al., CIKM 2009]

Important keywords as query    [Yang et al., WSDM 2009]
[Bendersky and Croft, WSDM 2009]

What query to formulate from important keywords?

information retrieval     text reuse

detection system     web search     query formulation

capacity constrained     search engine



## Candidate Document Retrieval for Web-Scale Text Reuse Detection*

Matthias Hagen and Benno Stein

Faculty of Media
Bauhaus-Universität Weimar, Germany
<first name>.<last name>@uni-weimar.de

**Abstract** Given a document d, the task of text reuse detection is to find those passages in d which in identical or paraphrased form also appear in other documents. To solve this problem at web-scale, keywords representing d's topics have to be combined to web queries. The retrieved web documents can then be delivered to a text reuse detection system for an in-depth analysis. We focus on the query formulation problem as the crucial first step in the detection process and present a new query formulation strategy that achieves convincing results: compared to a maximal termset query formulation strategy [10, 14], which is the most sensible non-heuristic baseline, we save on average 70% of the queries in realistic experiments. With respect to the candidate documents' quality, our heuristic retrieves documents that are, on average, more similar to the given document than the results of previously published query formulation strategies [4, 8].

information retrieval     text reuse

detection system     web search     query formulation

capacity constrained     search engine

**Candidate Document Retrieval for Web-Scale Text Reuse Detection***

Matthias Hagen and Benno Stein

Faculty of Media
Bauhaus-Universität Weimar, Germany
<first name>.<last name>@uni-weimar.de

**Abstract** Given a document *d*, the task of text reuse detection is to find those passages in *d* which in identical or paraphrased form also appear in other documents. To solve this problem at web-scale, keywords representing *d*'s topics have to be combined to web queries. The retrieved web documents can then be delivered to a text reuse detection system for an in-depth analysis. We focus on the query formulation problem as the crucial first step in the detection process and present a new query formulation strategy that achieves convincing results: compared to a maximal termset query formulation strategy [10, 14], which is the most sensible non-heuristic baseline, we save on average 70% of the queries in realistic experiments. With respect to the candidate documents' quality, our heuristic retrieves documents that are, on average, more similar to the given document than the results of previously published query formulation strategies [4, 8].

**Capacity-constrained Query Formulation**

Matthias Hagen and Benno Maria Stein

Faculty of Media
Bauhaus University Weimar, Germany

**Abstract** Given a set of keyphrases, we analyze how Web queries with these phrases can be formed that, taken altogether, return a specified number of hits. The use case of this problem is a plagiarism detection system that searches the Web for potentially plagiarized passages in a given suspicious document. For the query formulation problem we develop a heuristic search strategy based on co-occurrence probabilities. Compared to the maximal termset strategy [3], which can be considered as the most sensible non-heuristic baseline, our expected savings are on average 50% when queries for 8 or 10 phrases are to be constructed.

**1 Introduction**

The problem considered in this paper appears as an important sub task of automatic text plagiarism detection. Plagiarized passages in a suspicious document can be found via direct comparisons against potential source documents. Todays typical source of plagiarism is the Web, which obviously contains too many documents for direct comparisons. The straight-forward solution is to extract keyphrases from the suspicious document and to retrieve a tractable number of documents containing these phrases. These documents are considered

# Single keyword queries?

`information retrieval`

~~text reuse~~

~~detection system~~   ~~web search~~   ~~query formulation~~

~~capacity constrained~~   ~~search engine~~

# Single keyword queries?

~~information retrieval~~  ~~text reuse~~

**detection system**  ~~Web search~~  ~~query formulation~~

~~capacity constraint~~  ~~search engine~~

~~information retrieval~~                    ~~text reuse~~

**`detection system`**    ~~web search~~    ~~query formulation~~

~~capacity constrained~~    ~~search engine~~

information retrieval    text reuse

detection system    web search    query formulation

capacity constrained    search engine

information retrieval     text reuse

detection system     web search     query formulation

capacity constrained     search engine

information retrieval

detection system        web search

search engine

`information retrieval`     ~~t/e/h/t//w/e/b/s/h~~

`detection system`     `web search`     ~~x/q/u/e/n/y//f/l/h/t/m/d/l/a/k/i/o/d~~

~~h/a/p/a/c/i/t/y//k/o/n/s/t/w/a/i/n/e/d~~     `search engine`

What query to formulate from the keywords?

Not just one query!

Not just one query!

But a set of queries!

Not just one query!

But a set of queries!

Remark: Each returning not too many results . . .

information retrieval     text reuse

~~retrieval system~~     web search     query formulation

~~capacity constrained~~     search engine

~~information/retrieval~~ text reuse

detection system ~~web/search~~ query formulation

~~capacity/constrained~~ search engine

Google

"text reuse" "detection system" "query formulation" "search engine"

Advanced search

Search     About 5 results (0.16 seconds)

Everything

Images

Maps

Videos

News

Shopping

Books

More

[PDF] Candidate Document Retrieval for Web-Scale **Text Reuse**
**Detection** ▾
www.uni-weimar.de/medien/webis/publications/.../stein_2011i.pdf
File Format: PDF/Adobe Acrobat - Quick View
by M Hagen
ered to a **text reuse** detection **system** for an in-depth analysis. ... pared to a maximal
termset **query formulation** strategy [10, 14], which is the most ... ing **text reuse**
candidates on the **search engine's** ranking algorithm; potential **text reuse** ...

String Processing and Information Retrieval: 18th International ... -
**Google Books Result**
books.google.com/books?isbn=364224582X...
Roberto Grossi, Fabrizio Silvestri, Fabrizio Sebastiani - 2011 - Computers - 442 pages
We focus on the **query formulation** problem as the crucial first step in the ... A **text**
**reuse** detection **system** aims at finding passages within a given document which, ... of
how to query a web **search engine** using the extracted keywords. ...

Comparing query logs and pseudo-relevance feedbackfor web-search ...
portal.acm.org/citation.cfm?id=1277931
by RW White - 2007 - Cited by 10 - Related articles
Subjects: **Query formulation**. Additional Classification: ..... We evaluate our proposed
method on a commercial **search engine** log data. ...... In this paper we present a spam
**detection system** that combines ... ...... applications such as summarization, document
provenance, detecting **text reuse** and novelty detection. ...

All results
Sites with images
Related searches
Timeline
More search tools

~~information retrieval~~        ~~text mining~~

~~detection system~~        web search        query formulation

capacity constrained        search engine

~~information retrieval~~    ~~text mining~~

~~detection system~~    web search    query formulation

capacity constrained    search engine

# The 3 queries together . . .

## Properties

- All keywords covered                                 (similarity)
- Not too many results ($\leq 1000$)            (capacity)
- Desired document among the results     (quality)

## Problem

How to automatically find such query sets?

# Problem statement

## Capacity Constrained Query Formulation

- Given:
  1. Set $W$ of keywords
  2. Query interface for a web search engine
  3. Upper bound $k$ on the number of desired results

- Find a family $\mathcal{Q} \subseteq 2^W$ of queries:
  - returning $\leq k$ results
  - covering all keywords from $W$.

## Optimization Problem!

Minimize the number of submitted web queries to find $\mathcal{Q}$.

# Problem statement

## Capacity Constrained Query Formulation

- Given:
  1. Set $W$ of keywords
  2. Query interface for a web search engine
  3. Upper bound $k$ on the number of desired results

- Find a family $\mathcal{Q} \subseteq 2^W$ of queries:
  - returning $\leq k$ results
  - covering all keywords from $W$.

## Optimization Problem!

Minimize the number of submitted web queries to find $\mathcal{Q}$.

$\{w1, w2, w3, w4, w5\}$

*underflowing*

$\{w1, w2, w3, w4\}$  $\{w1, w2, w3, w5\}$  $\{w1, w2, w4, w5\}$  $\{w1, w3, w4, w5\}$  $\{w2, w3, w4, w5\}$

$\{w1,w2,w3\}$ $\{w1,w2,w4\}$ $\{w1,w2,w5\}$ $\{w1,w3,w4\}$ $\{w1,w3,w5\}$ $\{w1,w4,w5\}$ $\{w2,w3,w4\}$ $\{w2,w3,w5\}$ $\{w2,w4,w5\}$ $\{w3,w4,w5\}$

$\{w1, w2\}$  $\{w1, w3\}$  $\{w1, w4\}$  $\{w1, w5\}$  $\{w2, w3\}$  $\{w2, w4\}$  $\{w2, w5\}$  $\{w3, w4\}$  $\{w3, w5\}$  $\{w4, w5\}$

$\{w1\}$  $\{w2\}$  $\{w3\}$  $\{w4\}$  $\{w5\}$

*overflowing*

$\{ \}$

# Minimal non-overflowing queries



{w1, w2, w3, w4, w5}

*underflowing*

{w1, w2, w3, w4}    {w1, w2, w3, w5}    {w1, w2, w4, w5}    {w1, w3, w4, w5}    {w2, w3, w4, w5}

{w1,w2,w3}  {w1,w2,w4}  {w1,w2,w5}  {w1,w3,w4}  {w1,w3,w5}  {w1,w4,w5}  {w2,w3,w4}  {w2,w3,w5}  {w2,w4,w5}  {w3,w4,w5}

{w1, w2}   {w1, w3}   {w1, w4}   {w1, w5}   {w2, w3}   {w2, w4}   {w2, w5}   {w3, w4}   {w3, w5}   {w4, w5}

{w1}    {w2}    {w3}    {w4}    {w5}

*overflowing*

{ }

The baseline algorithm

Apriori

{w1, w2, w3, w4, w5}

*underflowing*

{w1, w2, w3, w4}    {w1, w2, w3, w5}    {w1, w2, w4, w5}    {w1, w3, w4, w5}    {w2, w3, w4, w5}

{w1,w2,w3}  {w1,w2,w4}  {w1,w2,w5}  {w1,w3,w4}  {w1,w3,w5}  {w1,w4,w5}  {w2,w3,w4}  {w2,w3,w5}  {w2,w4,w5}  {w3,w4,w5}

{w1, w2}  {w1, w3}  {w1, w4}  {w1, w5}  {w2, w3}  {w2, w4}  {w2, w5}  {w3, w4}  {w3, w5}  {w4, w5}

{w1}  {w2}  {w3}  {w4}  {w5}

*overflowing*

{ }

{w1, w2, w3, w4, w5}

*underflowing*

{w1, w2, w3, w4}    {w1, w2, w3, w5}    {w1, w2, w4, w5}    {w1, w3, w4, w5}    {w2, w3, w4, w5}

{w1,w2,w3}  {w1,w2,w4}  {w1,w2,w5}  {w1,w3,w4}  {w1,w3,w5}  {w1,w4,w5}  {w2,w3,w4}  {w2,w3,w5}  {w2,w4,w5}  {w3,w4,w5}

{w1, w2}    {w1, w3}    {w1, w4}    {w1, w5}    {w2, w3}    {w2, w4}    {w2, w5}    {w3, w4}    {w3, w5}    {w4, w5}

{w1}    {w2}    {w3}    {w4}    {w5}

*overflowing*

{ }

*underflowing*

{w1, w2, w3, w4, w5}

{w1, w2, w3, w4}    {w1, w2, w3, w5}    {w1, w2, w4, w5}    {w1, w3, w4, w5}    {w2, w3, w4, w5}

{w1,w2,w3} {w1,w2,w4} {w1,w2,w5} {w1,w3,w4} {w1,w3,w5} {w1,w4,w5} {w2,w3,w4} {w2,w3,w5} {w2,w4,w5} {w3,w4,w5}

{w1, w2}    {w1, w3}    {w1, w4}    {w1, w5}    {w2, w3}    {w2, w4}    {w2, w5}    {w3, w4}    {w3, w5}    {w4, w5}

{w1}    {w2}    {w3}    {w4}    {w5}

*overflowing*

{ }

# Baseline's Analysis

## Major drawback

All intermediate queries submitted. $\rightarrow$ Bad run time!

## Major drawback

All intermediate queries submitted. $\rightarrow$ Bad run time!

## Idea

Estimate the result list length before query submission.

The improved heuristic

Apriori + estimation

# Co-occurrences for estimation

Estimate: `"information retrieval"` `"query formulation"` + `"web search"`

# Co-occurrences for estimation

Estimate: `"information retrieval" "query formulation"` + `"web search"`

Known: `"information retrieval" "query formulation"`     87 100 results

# Co-occurrences for estimation

Estimate: `"information retrieval"` `"query formulation"` `+ "web search"`

Known:     `"information retrieval"` `"query formulation"`        87 100 results

              `"information retrieval"`        `+ "web search"`        16 % remain

# Co-occurrences for estimation

Estimate: `"information retrieval"` <span style="color:red">`"query formulation" + "web search"`</span>

Known:  `"information retrieval" "query formulation"`  87 100 results

`"information retrieval"    + "web search"`  16 % remain

`"query formulation"    + "web search"`  22 % remain

# Co-occurrences for estimation

Estimate: `"information retrieval" "query formulation" + "web search"`

Known: 

| | | |
|---|---|---|
| `"information retrieval" "query formulation"` | | 87 100 results |
| `"information retrieval"` | `+ "web search"` | 16 % remain |
| `"query formulation"` | `+ "web search"` | 22 % remain |

Our estimation scheme:

$$\mathrm{avg}(16\,\%\,,\,22\,\%) \;=\; 19\,\%$$
$$87\,100 \,\cdot\, 0.19 \;=\; 16\,500 \text{ results}$$

## Co-occurrences for estimation

Estimate: `"information retrieval" "query formulation" + "web search"`

Known:

| | |
|---|---|
| `"information retrieval" "query formulation"` | 87 100 results |
| `"information retrieval"`     `+ "web search"` | 16 % remain |
| `"query formulation"`     `+ "web search"` | 22 % remain |

Our estimation scheme:

$$\mathrm{avg}(16\,\%\,,\,22\,\%) \;=\; 19\,\%$$

$$87\,100 \cdot 0.19 \;=\; 16\,500 \text{ results}$$

Control:      Google     35 700 results

## Co-occurrences for estimation

Estimate: `"information retrieval" "query formulation" + "web search"`

Known: 

| | | |
|---|---|---|
| `"information retrieval"` `"query formulation"` | | 87 100 results |
| `"information retrieval"` `+ "web search"` | | 16 % remain |
| `"query formulation"` `+ "web search"` | | 22 % remain |

Our estimation scheme:

$$\mathrm{avg}(16\,\%\,,\,22\,\%) \;=\; 19\,\%$$
$$87\,100 \cdot 0.19 \;=\; 16\,500 \text{ results}$$

Control: Google 35 700 results

### Observation

Our scheme usually underestimates the real result list length.

What about performance?

# Experimental setup

## Corpus

- 257 pairs of two versions of papers
- 10 keywords from more mature version

## System

- Bing API as search engine
- Set $k = 1000$

# Experimental setup

## Corpus

- 257 pairs of two versions of papers
- 10 keywords from more mature version

## System

- Bing API as search engine
- Set $k = 1000$

# Baseline vs. heuristic

| Number of keywords | | 4 | 6 | 8 | 10 |
|---|---|---|---|---|---|
| complete query overflows | | 207 | 146 | 102 | 81 |
| $\mathcal{Q}$ computation possible | | 50 | 111 | 155 | 176 |
| Avg. queries submitted | heuristic | **6.69** | **13.30** | **32.58** | **95.86** |
| | baseline | 10.65 | 34.60 | 106.19 | 302.87 |

# Baseline vs. heuristic

What about the candidate document quality?

# Candidates' similarity to original document

| | **Approach** | | | |
| | Heuristic | Frequent | Rare | Random |
|---|---|---|---|---|
| 10 most similar doc's | 0.55 | 0.55 | **0.56** | **0.56** |
| 100 most similar doc's | **0.39** | 0.37 | 0.35 | 0.29 |
| all retrieved doc's | **0.29** | 0.25 | 0.22 | 0.21 |

Almost the end: The take-away messages!

## Results

- Candidate document retrieval
  - not just one query
  - set of queries $\rightarrow$ capacity
- Co-occurrence informed heuristic
- Good quality candidates

## Future work

- Which approach actually finds more text reuse?

## Results

- Candidate document retrieval
  - not just one query
  - set of queries $\rightarrow$ capacity
- Co-occurrence informed heuristic
- Good quality candidates

## Future work

- Which approach actually finds more text reuse?

# What we have (not) done

## Results

- Candidate document retrieval
  - not just one query
  - set of queries $\rightarrow$ capacity
- Co-occurrence informed heuristic
- Good quality candidates

## Future work

- Which approach actually finds more text reuse?

# Thank you
🙂