# Automatic Text Segmentation: TextTiling

Will Roberts
wroberts@coli.uni-sb.de

Tuesday, 20 November, 2007

**Introduction**
TextTiling
Evaluation
Conclusion

Motivation
Approaches

**Introduction**
TextTiling
Evaluation
Conclusion

**Motivation**
Approaches

## Uses for Automatic Text Segmentation

- hypertext display
- information/passage retrieval
- text summarization
- automatic text generation
- measuring stylistic variation for genre detection
- aligning parallel multilingual corpora
- breaking up connected documents

**Introduction**
TextTiling
Evaluation
Conclusion

Motivation
**Approaches**

## Other Text Segmentation Approaches

- Clustering or similarity matrices based on word co-occurance
- Machine-learning or hand-crafted solutions for detection of cue words

Introduction
**TextTiling**
Evaluation
Conclusion

**Theory**
Method

# TextTiling Theory

- Focus on multi-paragraph units in expository text
    - Topics not always contained in single paragraphs
- Identify subtopic shifts
    - Subtopic: piece of text "about" something
    - Identify topic shift, not topic
    - Linear segmentation
- Subtopic shifts associated with change in vocabulary
    - Linguistically simple: no prosody, discourse markers, pronoun reference resolution, . . .

Introduction
**TextTiling**
Evaluation
Conclusion

**Theory**
Method

# TextTiling Theory

- Focus on multi-paragraph units in expository text
    - Topics not always contained in single paragraphs
- Identify subtopic shifts
    - Subtopic: piece of text "about" something
    - Identify topic shift, not topic
    - Linear segmentation
- Subtopic shifts associated with change in vocabulary
    - Linguistically simple: no prosody, discourse markers, pronoun reference resolution, . . .

Introduction
**TextTiling**
Evaluation
Conclusion

**Theory**
Method

# TextTiling Theory

- Focus on multi-paragraph units in expository text
  - Topics not always contained in single paragraphs
- Identify subtopic shifts
  - Subtopic: piece of text "about" something
  - Identify topic shift, not topic
  - Linear segmentation
- Subtopic shifts associated with change in vocabulary
  - Linguistically simple: no prosody, discourse markers, pronoun reference resolution, . . .

Introduction
**TextTiling**
Evaluation
Conclusion

**Theory**
Method

# Word Occurance Counts

Introduction
**TextTiling**
Evaluation
Conclusion

Theory
**Method**

# TextTiling Algorithm

1. Tokenize
2. Calculate lexical similarity scores
3. Determine inter-sentence boundaries

Introduction
**TextTiling**
Evaluation
Conclusion

Theory
**Method**

## Lexical Similarity

- Computed for each sentence gap in text
- Measure of the lexical similarity of the two sentences/blocks on either side
- Moving window of size $k$

Introduction
**TextTiling**
Evaluation
Conclusion

Theory
**Method**

## Lexical Similarity

- Block comparison
  - Normalized inner product of two word vectors
- New vocabulary
  - New words introduced in a window centered around the sentence boundary
- Lexical chaining
  - Number of lexical chains active at the sentence boundary

Introduction
**Text Tiling**
Evaluation
Conclusion

Theory
**Method**

## Lexical Similarity

- Block comparison
    - Normalized inner product of two word vectors
- New vocabulary
    - New words introduced in a window centered around the sentence boundary
- Lexical chaining
    - Number of lexical chains active at the sentence boundary

Introduction
**TextTiling**
Evaluation
Conclusion

Theory
**Method**

## Lexical Similarity

- Block comparison
    - Normalized inner product of two word vectors
- New vocabulary
    - New words introduced in a window centered around the sentence boundary
- Lexical chaining
    - Number of lexical chains active at the sentence boundary

Introduction
**TextTiling**
Evaluation
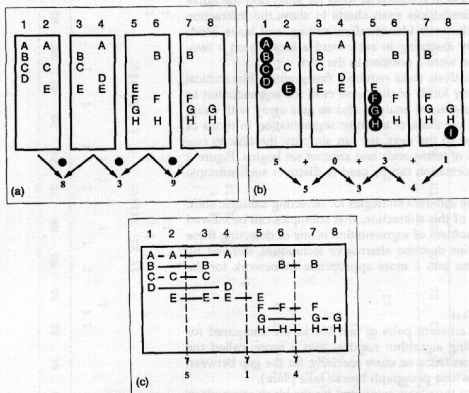Conclusion

Theory
**Method**

## Lexical Similarity



**Figure 3**
Illustration of three ways to compute the lexical score at gaps between sentences. Numbers indicate a numbered sequence of sentences, columns of letters signify the terms in the given sentence. (a) Blocks – dot product of vectors of word counts in the block on the left and the block on the right. (b) Vocabulary introduction – the number of words that occur for the first time within the interval centered at the sentence gap. (c) Chains – the number of active chains, or terms that repeat within threshold sentences and span the sentence gap.

Introduction
**TextTiling**
Evaluation
Conclusion

Theory
**Method**

# Calculating Word Significance Values

### Example

Saarland University (German: Universität des Saarlandes) is a university located in Saarbrücken, the capital of the German state of Saarland. It was founded in 1948 in co-operation with France and is organized in 8 faculties that cover all major fields of science. The university is particularly well known for research and education in Computer Science and Medicine.
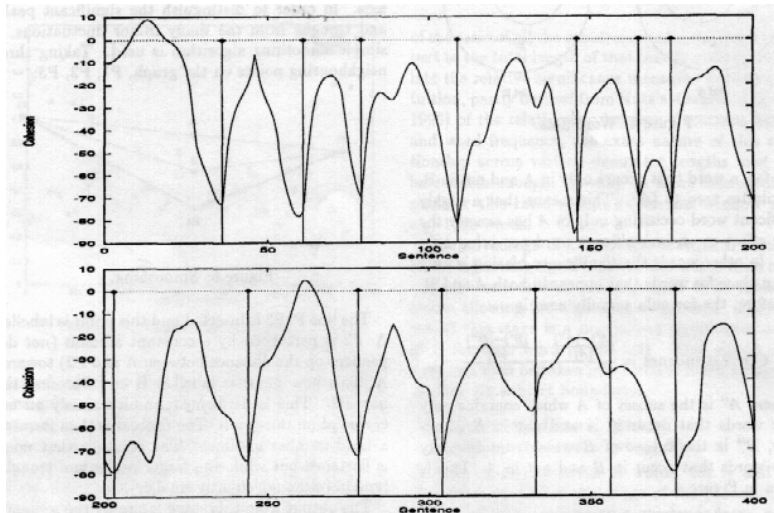
Saarland University, the first to be established after the Second World War, was founded in November 1948 with the support of the French Government and under the auspices of the University of Nancy.

At the time the Saarland found itself in the special situation of being partly autonomous and linked to France by economic . . .

Introduction
**TextTiling**
Evaluation
Conclusion

Theory
**Method**

## Boundary Identification

- *Depth score* is computed at each sentence gap from the lexical scores
    - Score is the sum of the heights of the peaks on either side of the sentence gap
    - $D_i = (s_{i-1} - s_i) + (s_{i+1} - s_i) = s_{i-1} + s_{i+1} - 2s_i$
- Smooth depth scores and choose local minima

Introduction
**TextTiling**
Evaluation
Conclusion

Theory
**Method**

# Boundary Identification

## Evaluation Caveats

- Algorithm depends on parameters: window size, smoothing, number of boundaries
- Evaluation depends on desired application: precision, recall, and *near-misses*
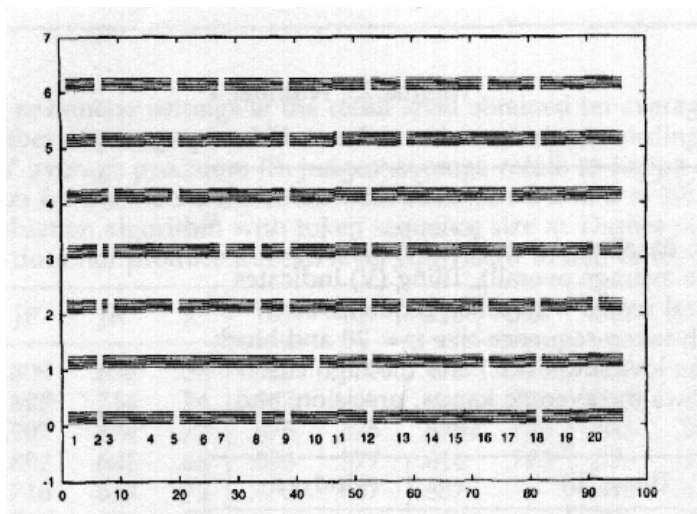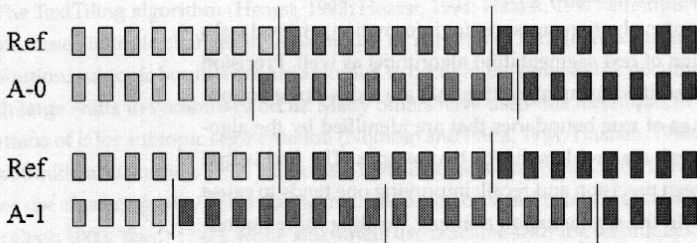
## Evaluation Methods

- Breaking consecutive documents
- Comparison with human judges
  - Humans rarely agree on correct segmentation
  - Consensus of human judges can be used as "gold standard"
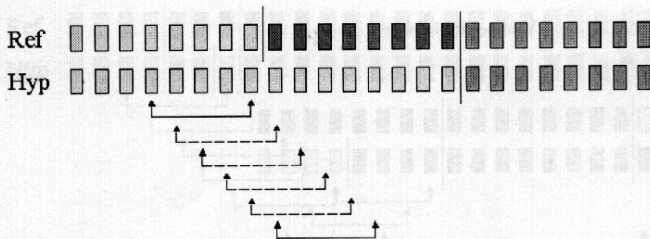- $P_k$ metric

## Evaluation Methods

- Breaking consecutive documents
- Comparison with human judges
    - Humans rarely agree on correct segmentation
    - Consensus of human judges can be used as "gold standard"
- $P_k$ metric

## Evaluation Methods

- Breaking consecutive documents
- Comparison with human judges
  - Humans rarely agree on correct segmentation
  - Consensus of human judges can be used as "gold standard"
- $P_k$ metric

# Human Judges

# $P_k$ Metric



**Figure 1**
Two hypothetical segmentations of the same reference (ground truth) document segmentation. The boxes indicate sentences or other units of subdivision, and spaces between boxes indicate potential boundary locations. Algorithm A-0 makes two near-misses, while Algorithm A-1 misses both boundaries by a wide margin and introduces three false positives. Both algorithms would receive scores of 0 for both precision and recall.

# $P_k$ Metric



**Figure 2**
An illustration of how the $P_k$ metric handles false negatives. The arrowed lines indicate the two poles of the probe as it moves from left to right, the boxes indicate sentences or other units of subdivision, and the width of the window ($k$) is four, meaning four potential boundaries fall between the two ends of the probe. Solid lines indicate no penalty is assigned, dashed lines indicate a penalty is assigned. Total penalty is always $k$ for false negatives.

## Conclusion

### Strengths

- Linguistically and computationally simple
- Language independent

### Weaknesses

- Designed for expository text; poor for narrative texts and discourse
- *Near-miss* errors might be unacceptable for some applications
- Cannot extract hierarchical text structure

## Literature

📄 Hearst, M.

Segmenting Text into Multi-paragraph Subtopic Passages

*Computational Linguistics*, 23(1), 1997.

📄 Pevzner, L., and Hearst, M.

A Critique and Improvement of an Evaluation Metric for Text Segmentation

*Computational Linguistics*, pp. 9–16, 2001.

📄 Richmond, K., Smith, A., and Amitay, E.

Detecting Subject Boundaries Within a Text: A Language Independent Statistical Approach

*EMNLP*, 1997.