

# How to Teach the Concepts of „Information Retrieval“ using Shell Commands

Andreas Schmidt

Institute for Automation and Applied Informatics  
Karlsruhe Institute of Technology  
Karlsruhe, Germany

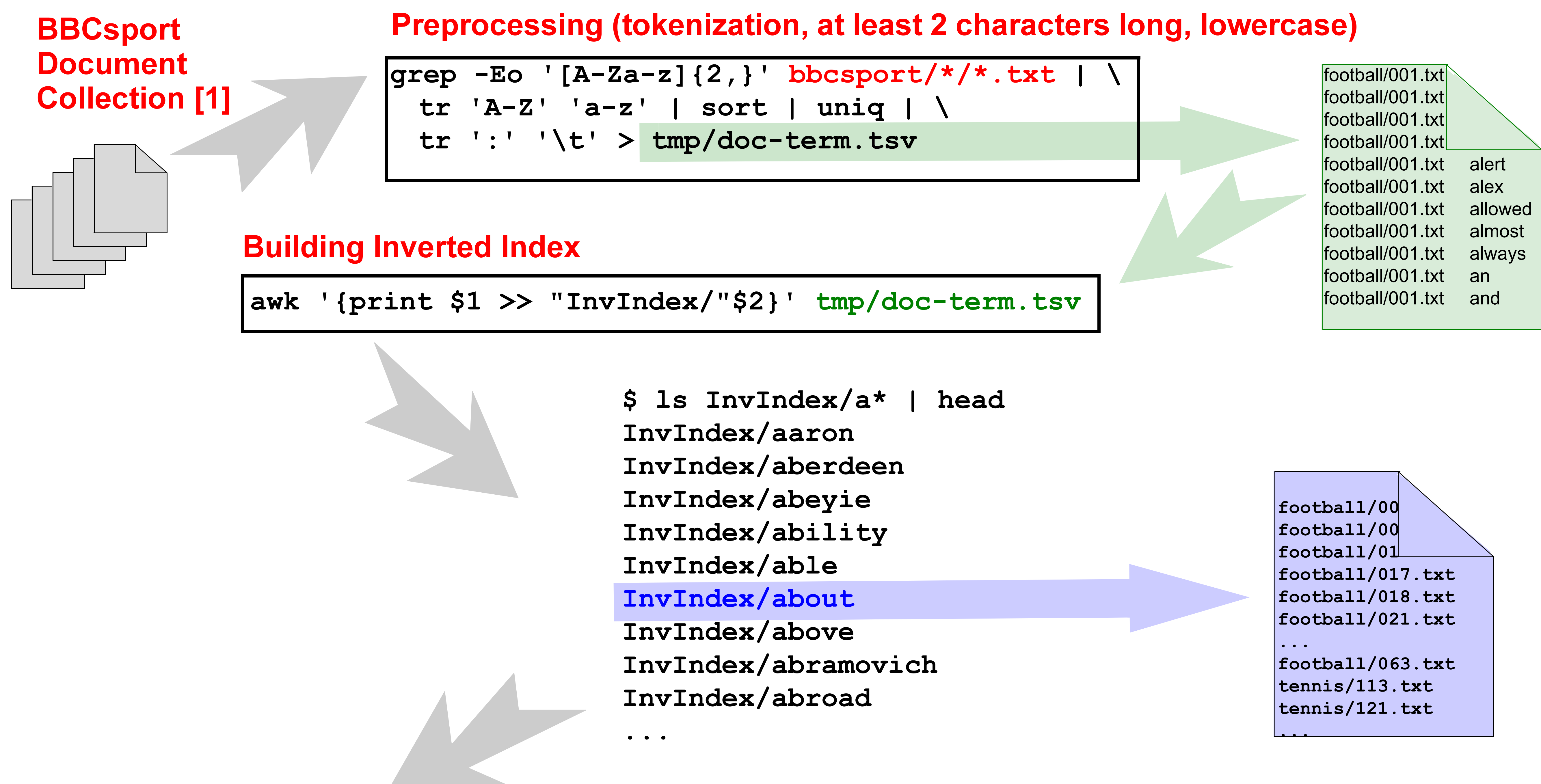
andreas.schmidt@kit.edu

Department of Computer Science and Business Information Systems  
Hochschule Karlsruhe  
University of Applied Sciences  
Karlsruhe, Germany  
andreas.schmidt@h-ka.de

## General Approach

- Using shell commands to get informations about the vocabulary (size of vocabulary for different languages, frequency of words, Zipfs Law)
  - Building a stopwordlist from a corpus
  - Construction of an Inverted Index based on ASCII files using shell commands
    - Boolean Retrieval
    - Ranking (using tf, idf)
    - Phrase match
  - Implementation of character-filter, tokenizer, token-filter using shell commands
  - Similarity of documents
  - Simple Entity Recognition
- Used tools:
    - grep
    - sed
    - cut
    - tr
    - uniq
    - sort
    - awk
    - cat
    - comm
    - join
    - less
    - wc
    - merge
    - split/csplit
    - head/tail
    - find
    - xargs

## Principle Idea of Index Construction



**Query: Looking for documents containing „ronaldo“ and „barcelona“:**

```
$ cat InvIndex/ronaldo InvIndex/barcelona | sort | uniq -d
```

```
bbc sport/football/057.txt
bbc sport/football/101.txt
...
```

documents containing both terms

[1] D. Greene and P. Cunningham. "Practical Solutions to the Problem of Diagonal Dominance in Kernel Document Clustering", Proc. ICML 2006.