

# Segmentation Of Layout-Based Documents

## Bachelor's Thesis

Albert-Ludwigs-Universität Freiburg



UNI  
FREIBURG

Elias Kempf

Department of Computer Science

Chair of Algorithms and Data Structures

October 20, 2021

# What problem do we want to solve?



# What problem do we want to solve?



- Extracting text blocks from PDFs

# What problem do we want to solve?



- Extracting text blocks from PDFs

## A Benchmark and Evaluation for Text Extraction from PDF

Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

### ABSTRACT

Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

### 1.1 Kinds of semantic information

In the following, we briefly describe the kind of semantic information that we investigate in this paper.

**Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

# What problem do we want to solve?



- Extracting text blocks from PDFs

## A Benchmark and Evaluation for Text Extraction from PDF

Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

### ABSTRACT

Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

### 1.1 Kinds of semantic information

In the following, we briefly describe the kind of semantic information that we investigate in this paper.

**Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

# What problem do we want to solve?



- Extracting text blocks from PDFs
- Sort extracted text blocks by natural reading order

## **A Benchmark and Evaluation for Text Extraction from PDF**

Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

### **ABSTRACT**

Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

### **1.1 Kinds of semantic information**

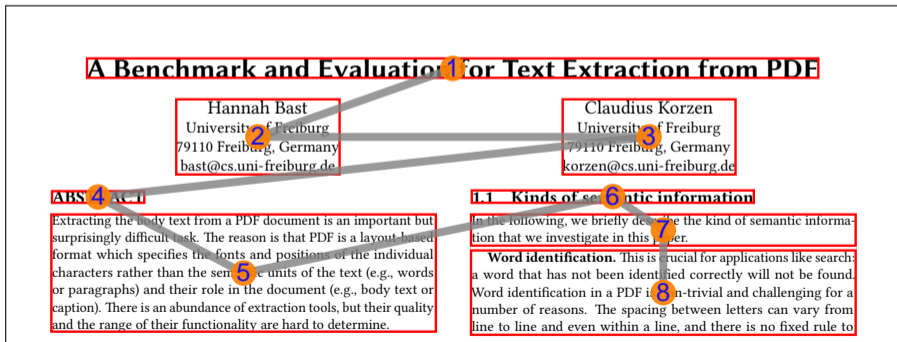
In the following, we briefly describe the kind of semantic information that we investigate in this paper.

**Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

# What problem do we want to solve?



- Extracting text blocks from PDFs
- Sort extracted text blocks by natural reading order



# Why is this problem difficult?





# Why is this problem difficult?



- PDF is a layout-based format (meaning it does not store text as words, lines, or paragraphs)

# Why is this problem difficult?



- PDF is a layout-based format (meaning it does not store text as words, lines, or paragraphs)
- Only characters, their bounding boxes, and font information is stored

# Why is this problem difficult?



- PDF is a layout-based format (meaning it does not store text as words, lines, or paragraphs)
- Only characters, their bounding boxes, and font information is stored
- Usually also no whitespace characters

# Why is this problem difficult?



## A Benchmark and Evaluation for Text Extraction from PDF

Hannah Bast  
University of Freiburg

Claudius Korzen  
University of Freiburg

# Why is this problem difficult?



## **A** Benchmark and Evaluation for Text Extraction from PDF

Hannah Bast  
University of Freiburg

Claudius Korzen  
University of Freiburg

Character	bounding box	font name	font size
„A“	(75.8, 697.2), (87.9, 708.5)	Arial	17.2

# Why is this problem difficult?



## A B Benchmark and Evaluation for Text Extraction from PDF

Hannah Bast  
University of Freiburg

Claudius Korzen  
University of Freiburg

Character	bounding box	font name	font size
„A“	(75.8, 697.2), (87.9, 708.5)	Arial	17.2
„B“	(92.2, 697.2), (103.8, 708.4)	Arial	17.2

# Why is this problem difficult?



## A Benchmark and Evaluation for Text Extraction from PDF

Hannah Bast  
University of Freiburg

Claudius Korzen  
University of Freiburg

Character	bounding box	font name	font size
„A“	(75.8, 697.2), (87.9, 708.5)	Arial	17.2
„B“	(92.2, 697.2), (103.8, 708.4)	Arial	17.2
„e“	(103.8, 697.1), (112.5, 704.8)	Arial	17.2

# Why is this problem difficult?



## A Benchmark and Evaluation for Text Extraction from PDF

Hannah Bast  
University of Freiburg

Claudius Korzen  
University of Freiburg

Character	bounding box	font name	font size
„A“	(75.8, 697.2), (87.9, 708.5)	Arial	17.2
„B“	(92.2, 697.2), (103.8, 708.4)	Arial	17.2
„e“	(103.8, 697.1), (112.5, 704.8)	Arial	17.2
	⋮		



# Why is this problem difficult?



- Reading order is also often difficult to detect

# Why is this problem difficult?



- Reading order is also often difficult to detect
- Especially, in documents featuring a two-column layout:

## **A Benchmark and Evaluation for Text Extraction from PDF**

Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

### **ABSTRACT**

Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

### **1.1 Kinds of semantic information**

In the following, we briefly describe the kind of semantic information that we investigate in this paper.

**Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to



For each page of a given PDF, ...

- Our input is a list of characters. Each character comes with its bounding box and its font information.
- Our output is a list of text blocks sorted by reading order.

# Approach outline



**UNI  
FREIBURG**

# Approach outline



- We separate our approach into two main steps

# Approach outline



- We separate our approach into two main steps
- First, we use page segmentation to detect text blocks

# Approach outline



- We separate our approach into two main steps
- First, we use page segmentation to detect text blocks
- Second, we order the detected text blocks using a similar but more informed approach



- Page segmentation is the process of reassembling the characters of a layout-based document into semantic units like words, lines, or text blocks





- Page segmentation is the process of reassembling the characters of a layout-based document into semantic units like words, lines, or text blocks
- We only focus on reassembling characters into text blocks



- Page segmentation is the process of reassembling the characters of a layout-based document into semantic units like words, lines, or text blocks
- We only focus on reassembling characters into text blocks
- We perform our segmentation using an XY-cut algorithm

# XY-cut algorithm



- An XY-cut algorithm can be used to group the characters of a page



- An XY-cut algorithm can be used to group the characters of a page
- It does so by applying vertical cuts (through the **X**-axis) and horizontal cuts (through the **Y**-axis) to the page

- An XY-cut algorithm can be used to group the characters of a page
- It does so by applying vertical cuts (through the **X**-axis) and horizontal cuts (through the **Y**-axis) to the page
- Diagonal cuts and cuts through characters are not allowed



- An XY-cut algorithm can be used to group the characters of a page
- It does so by applying vertical cuts (through the **X**-axis) and horizontal cuts (through the **Y**-axis) to the page
- Diagonal cuts and cuts through characters are not allowed
- Cuts can also be used to detect reading order (more on that later)

## A Benchmark and Evaluation for Text Extraction from PDF

Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

### ABSTRACT

Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

### 1.1 Kinds of semantic information

In the following, we briefly describe the kind of semantic information that we investigate in this paper.

**Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

## A Benchmark and Evaluation for Text Extraction from PDF

Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

### ABSTRACT

Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

### 1.1 Kinds of semantic information

In the following, we briefly describe the kind of semantic information that we investigate in this paper.

**Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to



## A Benchmark and Evaluation for Text Extraction from PDF

Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

### ABSTRACT

Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

### 1.1 Kinds of semantic information

In the following, we briefly describe the kind of semantic information that we investigate in this paper.

**Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

## A Benchmark and Evaluation for Text Extraction from PDF

Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

### ABSTRACT

Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

### 1.1 Kinds of semantic information

In the following, we briefly describe the kind of semantic information that we investigate in this paper.

**Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

## A Benchmark and Evaluation for Text Extraction from PDF

Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

### ABSTRACT

Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

### 1.1 Kinds of semantic information

In the following, we briefly describe the kind of semantic information that we investigate in this paper.

**Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

## A Benchmark and Evaluation for Text Extraction from PDF

Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

### ABSTRACT

Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

### 1.1 Kinds of semantic information

In the following, we briefly describe the kind of semantic information that we investigate in this paper.

**Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

## A Benchmark and Evaluation for Text Extraction from PDF

Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

### ABSTRACT

Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

### 1.1 Kinds of semantic information

In the following, we briefly describe the kind of semantic information that we investigate in this paper.

**Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

## A Benchmark and Evaluation for Text Extraction from PDF

Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

### ABSTRACT

Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

### 1.1 Kinds of semantic information

In the following, we briefly describe the kind of semantic information that we investigate in this paper.

**Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

## A Benchmark and Evaluation for Text Extraction from PDF

Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

### ABSTRACT

Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

### 1.1 Kinds of semantic information

In the following, we briefly describe the kind of semantic information that we investigate in this paper.

**Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

# Reading order detection



- Each cut divides a page into two parts



# Reading order detection



- Each cut divides a page into two parts
- We order them respecting our top-to-bottom left-to-right writing system

- Each cut divides a page into two parts
- We order them respecting our top-to-bottom left-to-right writing system

## A Benchmark and Evaluation for Text Extraction from PDF

Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

### ABSTRACT

Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

### 1.1 Kinds of semantic information

In the following, we briefly describe the kind of semantic information that we investigate in this paper.

**Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

# Reading order detection



- Detecting text blocks also yields a preliminary reading order

# Reading order detection



- Detecting text blocks also yields a preliminary reading order
- Problem: cuts are not specifically chosen for reading order detection

- Detecting text blocks also yields a preliminary reading order
- Problem: cuts are not specifically chosen for reading order detection

## A Benchmark and Evaluation for Text Extraction from PDF

Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

### ABSTRACT

Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

### 1.1 Kinds of semantic information

In the following, we briefly describe the kind of semantic information that we investigate in this paper.

**Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

- Detecting text blocks also yields a preliminary reading order
- Problem: cuts are not specifically chosen for reading order detection

## A Benchmark and Evaluation for Text Extraction from PDF

Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

### ABSTRACT

Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

### 1.1 Kinds of semantic information

In the following, we briefly describe the kind of semantic information that we investigate in this paper.

**Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

# Reading order detection



- Solution: detect reading order **after** detecting text blocks

- Solution: detect reading order **after** detecting text blocks
- This allows us to make use of information about the text blocks themselves:

**A Benchmark and Evaluation for Text Extraction from PDF**

Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

**ABSTRACT**  
Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

**1.1 Kinds of semantic information**  
In the following, we briefly describe the kind of semantic information that we investigate in this paper.  
**Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to



- Solution: detect reading order **after** detecting text blocks
- This allows us to make use of information about the text blocks themselves:

	title	<b>A Benchmark and Evaluation for Text Extraction from PDF</b>			
	author	Hannah Bast University of Freiburg 79110 Freiburg, Germany bast@cs.uni-freiburg.de	author	Claudius Korzen University of Freiburg 79110 Freiburg, Germany korzen@cs.uni-freiburg.de	
para- graph	<b>ABSTRACT:</b> heading	Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.	<b>1.1 Kinds of semantic information</b> heading	in the following, we briefly describe the kind of semantic information that we investigate in this paper.	para- graph
			<b>Word identification.</b> This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to	para- graph	

- Solution: detect reading order **after** detecting text blocks
- This allows us to make use of information about the text blocks themselves:

title	
<b>A Benchmark and Evaluation for Text Extraction from PDF</b>	
author	Hannah Bast University of Freiburg 79110 Freiburg, Germany bast@cs.uni-freiburg.de
author	Claudius Korzen University of Freiburg 79110 Freiburg, Germany korzen@cs.uni-freiburg.de
para-graph	<b>ABSTRACT:</b> heading Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.
	<b>1.1 Kinds of semantic information</b> heading In the following, we briefly describe the kind of semantic information that we investigate in this paper. <b>Word identification.</b> This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to
	para-graph para-graph

- Solution: detect reading order **after** detecting text blocks
- This allows us to make use of information about the text blocks themselves:

title	
<b>A Benchmark and Evaluation for Text Extraction from PDF</b>	
author	Hannah Bast University of Freiburg 79110 Freiburg, Germany bast@cs.uni-freiburg.de
author	Claudius Korzen University of Freiburg 79110 Freiburg, Germany korzen@cs.uni-freiburg.de
para-graph	<b>ABSTRACT</b> heading Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.
	<b>1.1 Kinds of semantic information</b> heading In the following, we briefly describe the kind of semantic information that we investigate in this paper. <b>Word identification.</b> This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to
	para-graph para-graph



- We predict the semantic roles of detected text blocks using a machine-learning model developed by Korzen



- We predict the semantic roles of detected text blocks using a machine-learning model developed by Korzen
- We then use the XY-cut algorithm again but now choose cuts using a machine-learning model which also considers semantic roles



- We predict the semantic roles of detected text blocks using a machine-learning model developed by Korzen
- We then use the XY-cut algorithm again but now choose cuts using a machine-learning model which also considers semantic roles
- This way, we can correct potential mistakes in the preliminary reading order

# Evaluation setup



**UNI  
FREIBURG**



- We evaluate our approach on 1,750 randomly selected articles from `arXiv.org`





- We evaluate our approach on 1,750 randomly selected articles from arXiv.org
- We evaluate both text block detection and reading order detection

- We evaluate our approach on 1,750 randomly selected articles from `arXiv.org`
- We evaluate both text block detection and reading order detection
- We compute the expected text blocks and their natural reading order from  $\text{T}_{\text{E}}\text{X}$  data

- We evaluate our approach on 1,750 randomly selected articles from arXiv.org
- We evaluate both text block detection and reading order detection
- We compute the expected text blocks and their natural reading order from  $\text{T}_{\text{E}}\text{X}$  data
- We compare text blocks using their bounding boxes



- Text block detection:



- Text block detection:
  - $B_G^=$  := percentage of expected blocks that were detected



- Text block detection:
  - $B_G^-$  := percentage of expected blocks that were detected
  - $B_A^-$  := percentage of detected blocks that were expected

- Text block detection:
  - $B_G^-$  := percentage of expected blocks that were detected
  - $B_A^-$  := percentage of detected blocks that were expected
  - $B_G^+$  := percentage of expected blocks that were split too much

- Text block detection:
  - $B_G^-$  := percentage of expected blocks that were detected
  - $B_A^-$  := percentage of detected blocks that were expected
  - $B_G^+$  := percentage of expected blocks that were split too much
  - $B_A^+$  := percentage of detected blocks that were split too less



- Text block detection:
  - $B_G^-$  := percentage of expected blocks that were detected
  - $B_A^-$  := percentage of detected blocks that were expected
  - $B_G^+$  := percentage of expected blocks that were split too much
  - $B_A^+$  := percentage of detected blocks that were split too less
  
- Reading order detection:

- Text block detection:
  - $B_G^-$  := percentage of expected blocks that were detected
  - $B_A^-$  := percentage of detected blocks that were expected
  - $B_G^+$  := percentage of expected blocks that were split too much
  - $B_A^+$  := percentage of detected blocks that were split too less
- Reading order detection:
  - $\tau_n$  := the normalized Kendall- $\tau$ -correlation between expected and detected reading order

# Kendall- $\tau$ -correlation



- $\tau$  can be used to compare the order of a sequence of numbers to an ascending order

# Kendall- $\tau$ -correlation



- $\tau$  can be used to compare the order of a sequence of numbers to an ascending order
- To compute  $\tau$ , we count concordant and discordant pairs in a given sequence

- $\tau$  can be used to compare the order of a sequence of numbers to an ascending order
- To compute  $\tau$ , we count concordant and discordant pairs in a given sequence
- Let's look at an example:

7, 5, 6, 9

- $\tau$  can be used to compare the order of a sequence of numbers to an ascending order
- To compute  $\tau$ , we count concordant and discordant pairs in a given sequence
- Let's look at an example:

7, 5, 6, 9

Concordant pairs: 0

Discordant pairs: 0

- $\tau$  can be used to compare the order of a sequence of numbers to an ascending order
- To compute  $\tau$ , we count concordant and discordant pairs in a given sequence
- Let's look at an example:

$\overbrace{7, 5}, 6, 9$

Concordant pairs: 0

Discordant pairs: 1

- $\tau$  can be used to compare the order of a sequence of numbers to an ascending order
- To compute  $\tau$ , we count concordant and discordant pairs in a given sequence
- Let's look at an example:

$\overbrace{7, 5, 6, 9}$

Concordant pairs: 0

Discordant pairs: 2



- $\tau$  can be used to compare the order of a sequence of numbers to an ascending order
- To compute  $\tau$ , we count concordant and discordant pairs in a given sequence
- Let's look at an example:

$\overbrace{7, 5, 6, 9}$

Concordant pairs: 1

Discordant pairs: 2

- $\tau$  can be used to compare the order of a sequence of numbers to an ascending order
- To compute  $\tau$ , we count concordant and discordant pairs in a given sequence
- Let's look at an example:

7,  $\overbrace{5, 6}^{\text{concordant}}$ , 9

Concordant pairs: 2

Discordant pairs: 2

- $\tau$  can be used to compare the order of a sequence of numbers to an ascending order
- To compute  $\tau$ , we count concordant and discordant pairs in a given sequence
- Let's look at an example:

7,  $\overbrace{5, 6, 9}$

Concordant pairs: 3

Discordant pairs: 2

- $\tau$  can be used to compare the order of a sequence of numbers to an ascending order
- To compute  $\tau$ , we count concordant and discordant pairs in a given sequence
- Let's look at an example:

7, 5,  $\overbrace{6, 9}$

Concordant pairs: 4

Discordant pairs: 2

- Assuming the sequence does not contain duplicates, we can define  $\tau$  as

$$\tau = \frac{\#con - \#dis}{\#con + \#dis}$$

- Assuming the sequence does not contain duplicates, we can define  $\tau$  as

$$\tau = \frac{\#con - \#dis}{\#con + \#dis}$$

- For our example, we obtained  $\#con = 4$  and  $\#dis = 2$  yielding

$$\tau = \frac{4 - 2}{4 + 2} = \frac{2}{6} = \frac{1}{3}$$

- Assuming the sequence does not contain duplicates, we can define  $\tau$  as

$$\tau = \frac{\#con - \#dis}{\#con + \#dis}$$

- For our example, we obtained  $\#con = 4$  and  $\#dis = 2$  yielding

$$\tau = \frac{4 - 2}{4 + 2} = \frac{2}{6} = \frac{1}{3}$$

- $\tau$  takes values between -1 and 1, so we normalize it using

$$\tau_n = \frac{\tau + 1}{2}$$

- Average metric values on our evaluation dataset:

$B_G^-$	$B_A^-$	$B_G^+$	$B_A^-$	$\tau_n$	$\tau_n^f$
51.4%	46.7%	12.9%	14.7%	0.873	0.994



```
1  {
2    "glyphs": [{
3      "char": "A",
4      "font size": "11pt",
5      "bounding box": [1, 4, 2.5, 6]
6    },
7    {
8      "char": "4",
9      "font size": "11pt",
10     "bounding box": [1, 1, 3, 3]
11   }]
12 }
```

- The example omitted some important aspects:
  - How do we compute potential cuts algorithmically?  
⇒ using projection profiles of bounding boxes
  - How do we decide which cut to choose?  
⇒ based on cut size
  - When do we stop cutting?  
⇒ after cut size falls below a certain size threshold

### 1.1 Kinds of semantic information

In the following, we briefly describe the kind of semantic information that we investigate in this paper.

**Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

### 1.1 Kinds of semantic information

In the following, we briefly describe the kind of semantic information that we investigate in this paper.

**Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

### 1.1 Kinds of semantic information

In the following, we briefly describe the kind of semantic information that we investigate in this paper.

**Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

## 1 A Benchmark and Evaluation for Text Extraction from PDF

2

Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

3

Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

### 4 ABSTRACT

5 Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

### 6 1.1 Kinds of semantic information

7 In the following, we briefly describe the kind of semantic information that we investigate in this paper.

8 **Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

## 1 A Benchmark and Evaluation for Text Extraction from PDF

2

4

3 Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

7

Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

### 5 ABSTRACT

6 Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

### 8 1.1 Kinds of semantic information

9 In the following, we briefly describe the kind of semantic information that we investigate in this paper.

9 **Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

1 **A Benchmark and Evaluation for Text Extraction from PDF**

2 Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

3 Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

4 **ABSTRACT**

5 Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

6 **1.1 Kinds of semantic information**

7 In the following, we briefly describe the kind of semantic information that we investigate in this paper.

8 **Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

1 **A Benchmark and Evaluation for Text Extraction from PDF**

2

3 Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

4

7 Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

5 **ABSTRACT**

6 Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

8 **1.1 Kinds of semantic information**

9 In the following, we briefly describe the kind of semantic information that we investigate in this paper.

9 **Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

■  $\#exp = 8$

1 **A Benchmark and Evaluation for Text Extraction from PDF**

2 Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

3 Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

4 **ABSTRACT**

5 Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

6 **1.1 Kinds of semantic information**

7 In the following, we briefly describe the kind of semantic information that we investigate in this paper.

8 **Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

1 **A Benchmark and Evaluation for Text Extraction from PDF**

2

3 Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

4

7 Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

5 **ABSTRACT**

6 Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

8 **1.1 Kinds of semantic information**

9 In the following, we briefly describe the kind of semantic information that we investigate in this paper.

9 **Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

■  $\#exp = 8, \#det = 9$



1 **A Benchmark and Evaluation for Text Extraction from PDF**

2 Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

3 Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

4 **ABSTRACT**

5 Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

6 **1.1 Kinds of semantic information**

7 In the following, we briefly describe the kind of semantic information that we investigate in this paper.

8 **Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

1 **A Benchmark and Evaluation for Text Extraction from PDF**

2

3 Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

4

7 Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

5 **ABSTRACT**

6 Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

8 **1.1 Kinds of semantic information**

9 In the following, we briefly describe the kind of semantic information that we investigate in this paper.

9 **Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

■  $\#exp = 8, \#det = 9, \#cor = 4$

1 **A Benchmark and Evaluation for Text Extraction from PDF**

2 Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

3 Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

4 **ABSTRACT**

5 Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

6 **1.1 Kinds of semantic information**

7 In the following, we briefly describe the kind of semantic information that we investigate in this paper.

8 **Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

1 **A Benchmark and Evaluation for Text Extraction from PDF**

2

3 Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

4

5 **ABSTRACT**

6 Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

7 Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

8 **1.1 Kinds of semantic information**

9 In the following, we briefly describe the kind of semantic information that we investigate in this paper.

9 **Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

■  $\#exp = 8, \#det = 9, \#cor = 4$

$$B_G^- = \frac{\#cor}{\#exp} = \frac{4}{8} = 0.5 \hat{=} 50\%$$

$$B_A^- = \frac{\#cor}{\#det} = \frac{4}{9} = 0.\bar{4} \hat{=} 44.\bar{4}\%$$

1 **A Benchmark and Evaluation for Text Extraction from PDF**

2 Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

3 Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

4 **ABSTRACT**

5 Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

6 **1.1 Kinds of semantic information**

7 In the following, we briefly describe the kind of semantic information that we investigate in this paper.

8 **Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

1 **A Benchmark and Evaluation for Text Extraction from PDF**

2

3 Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

4

5 **ABSTRACT**

6 Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

7 Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

8 **1.1 Kinds of semantic information**

9 In the following, we briefly describe the kind of semantic information that we investigate in this paper.

9 **Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

1 **A Benchmark and Evaluation for Text Extraction from PDF**

2 Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

3 Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

4 **ABSTRACT**

5 Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

6 **1.1 Kinds of semantic information**

7 In the following, we briefly describe the kind of semantic information that we investigate in this paper.

8 **Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

1 **A Benchmark and Evaluation for Text Extraction from PDF**

2

3 Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

4

5 **ABSTRACT**

6 Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

7 Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

8 **1.1 Kinds of semantic information**

9 In the following, we briefly describe the kind of semantic information that we investigate in this paper.

9 **Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

■  $\#exp = 8$

1 **A Benchmark and Evaluation for Text Extraction from PDF**

2 Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

3 Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

4 **ABSTRACT**

5 Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

6 **1.1 Kinds of semantic information**

7 In the following, we briefly describe the kind of semantic information that we investigate in this paper.

8 **Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

1 **A Benchmark and Evaluation for Text Extraction from PDF**

2

3 Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

4

5 **ABSTRACT**

6 Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

7 Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

8 **1.1 Kinds of semantic information**

9 In the following, we briefly describe the kind of semantic information that we investigate in this paper.

9 **Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

■  $\#exp = 8, \#det = 9$

1 **A Benchmark and Evaluation for Text Extraction from PDF**

2 Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

3 Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

4 **ABSTRACT**

5 Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

6 **1.1 Kinds of semantic information**

7 In the following, we briefly describe the kind of semantic information that we investigate in this paper.

8 **Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

1 **A Benchmark and Evaluation for Text Extraction from PDF**

2

3 Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

4

5 **ABSTRACT**

6 Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

7 Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

8 **1.1 Kinds of semantic information**

9 In the following, we briefly describe the kind of semantic information that we investigate in this paper.

9 **Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

■  $\#exp = 8, \#det = 9, \#stm = 2$

1 **A Benchmark and Evaluation for Text Extraction from PDF**

2 Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

3 Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

4 **ABSTRACT**

5 Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

6 **1.1 Kinds of semantic information**

7 In the following, we briefly describe the kind of semantic information that we investigate in this paper.

8 **Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

1 **A Benchmark and Evaluation for Text Extraction from PDF**

2

3 Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

4

7 Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

5 **ABSTRACT**

6 Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

8 **1.1 Kinds of semantic information**

9 In the following, we briefly describe the kind of semantic information that we investigate in this paper.

9 **Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

■  $\#exp = 8$ ,  $\#det = 9$ ,  $\#stm = 2$ ,  $\#stl = 1$

1 **A Benchmark and Evaluation for Text Extraction from PDF**

2 Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

3 Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

4 **ABSTRACT**

5 Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

6 **1.1 Kinds of semantic information**

7 In the following, we briefly describe the kind of semantic information that we investigate in this paper.

8 **Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

1 **A Benchmark and Evaluation for Text Extraction from PDF**

2

3 Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

4

5 **ABSTRACT**

6 Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

7 Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

8 **1.1 Kinds of semantic information**

9 In the following, we briefly describe the kind of semantic information that we investigate in this paper.

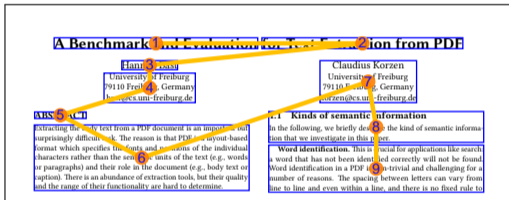
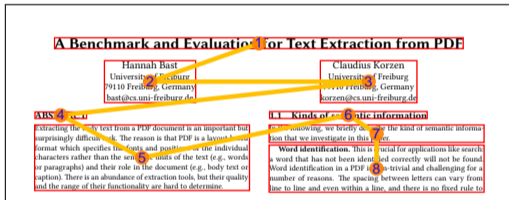
9 **Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

■  $\#exp = 8$ ,  $\#det = 9$ ,  $\#stm = 2$ ,  $\#stl = 1$

$$B_G^+ = \frac{\#stm}{\#exp} = \frac{2}{8} = 0.25 \hat{=} 25\%$$

$$B_A^- = \frac{\#stl}{\#det} = \frac{1}{9} = 0.\bar{1} \hat{=} 11.\bar{1}\%$$





## A Benchmark and Evaluation for Text Extraction from PDF

Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

ABSTRACT

Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

### 1.1 Kinds of semantic information

In the following, we briefly describe the kind of semantic information that we investigate in this paper.

**Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

## A Benchmark and Evaluation for Text Extraction from PDF

Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

ABSTRACT

Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

### 1.1 Kinds of semantic information

In the following, we briefly describe the kind of semantic information that we investigate in this paper.

**Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

## ■ Detected sequence:

## A Benchmark and Evaluation for Text Extraction from PDF

Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

ABSTRACT

Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

### 1.1 Kinds of semantic information

In the following, we briefly describe the kind of semantic information that we investigate in this paper.

**Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

## A Benchmark and Evaluation for Text Extraction from PDF

Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

ABSTRACT

Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

### 1.1 Kinds of semantic information

In the following, we briefly describe the kind of semantic information that we investigate in this paper.

**Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

- Detected sequence: 7

## A Benchmark and Evaluation for Text Extraction from PDF

Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

ABSTRACT

Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

### 1.1 Kinds of semantic information

In the following, we briefly describe the kind of semantic information that we investigate in this paper.

**Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

## A Benchmark and Evaluation for Text Extraction from PDF

Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

ABSTRACT

Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

### 1.1 Kinds of semantic information

In the following, we briefly describe the kind of semantic information that we investigate in this paper.

**Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

- Detected sequence: 7, 5

## A Benchmark and Evaluation for Text Extraction from PDF

Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

ABSTRACT

Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

### 1.1 Kinds of semantic information

In the following, we briefly describe the kind of semantic information that we investigate in this paper.

**Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

## A Benchmark and Evaluation for Text Extraction from PDF

Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

ABSTRACT

Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

### 1.1 Kinds of semantic information

In the following, we briefly describe the kind of semantic information that we investigate in this paper.

**Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

- Detected sequence: 7, 5, 6

## A Benchmark and Evaluation for Text Extraction from PDF

Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

ABSTRACT

Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

### 1.1 Kinds of semantic information

In the following, we briefly describe the kind of semantic information that we investigate in this paper.

**Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

## A Benchmark and Evaluation for Text Extraction from PDF

Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

ABSTRACT

Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

### 1.1 Kinds of semantic information

In the following, we briefly describe the kind of semantic information that we investigate in this paper.

**Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

- Detected sequence: 7, 5, 6, 9

## A Benchmark and Evaluation for Text Extraction from PDF

Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

ABSTRACT

Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

### 1.1 Kinds of semantic information

In the following, we briefly describe the kind of semantic information that we investigate in this paper.

**Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

## A Benchmark and Evaluation for Text Extraction from PDF

Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

ABSTRACT

Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

### 1.1 Kinds of semantic information

In the following, we briefly describe the kind of semantic information that we investigate in this paper.

**Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

- Detected sequence: 7, 5, 6, 9      #con = 4

## A Benchmark and Evaluation for Text Extraction from PDF

Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

ABSTRACT

Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

### 1.1 Kinds of semantic information

In the following, we briefly describe the kind of semantic information that we investigate in this paper.

**Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

## A Benchmark and Evaluation for Text Extraction from PDF

Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

ABSTRACT

Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

### 1.1 Kinds of semantic information

In the following, we briefly describe the kind of semantic information that we investigate in this paper.

**Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

- Detected sequence: 7, 5, 6, 9      #con = 4, #dis = 2



## A Benchmark and Evaluation for Text Extraction from PDF

Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

ABSTRACT

Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

### 1.1 Kinds of semantic information

In the following, we briefly describe the kind of semantic information that we investigate in this paper.

**Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

## A Benchmark and Evaluation for Text Extraction from PDF

Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

ABSTRACT

Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

### 1.1 Kinds of semantic information

In the following, we briefly describe the kind of semantic information that we investigate in this paper.

**Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

- Detected sequence: 7, 5, 6, 9       $\#con = 4, \#dis = 2$        $\tau = \frac{1}{3}$

**A Benchmark and Evaluation for Text Extraction from PDF**

Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

**ABS 4** Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

**1.1 Kinds of semantic information**  
In the following, we briefly describe the kind of semantic information that we investigate in this paper.

**Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

**A Benchmark and Evaluation for Text Extraction from PDF**

Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

**ABS 5** Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

**1.1 Kinds of semantic information**  
In the following, we briefly describe the kind of semantic information that we investigate in this paper.

**Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

■ Detected sequence: 7, 5, 6, 9      #con = 4, #dis = 2       $\tau = \frac{1}{3}$

■ We obtain

$$\tau_n = \frac{\tau + 1}{2}$$

## A Benchmark and Evaluation for Text Extraction from PDF

Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

ABS 4

Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

### 1.1 Kinds of semantic information

In the following, we briefly describe the kind of semantic information that we investigate in this paper.

**Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

## A Benchmark and Evaluation for Text Extraction from PDF

Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

ABS 5

Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

### 1.1 Kinds of semantic information

In the following, we briefly describe the kind of semantic information that we investigate in this paper.

**Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

- Detected sequence: 7, 5, 6, 9       $\#con = 4$ ,  $\#dis = 2$        $\tau = \frac{1}{3}$

- We obtain

$$\tau_n = \frac{\tau + 1}{2} = \frac{1/3 + 1}{2}$$

## A Benchmark and Evaluation for Text Extraction from PDF

Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

ABS 4

Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

### 1.1 Kinds of semantic information

In the following, we briefly describe the kind of semantic information that we investigate in this paper.

**Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

## A Benchmark and Evaluation for Text Extraction from PDF

Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

ABS 5

Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

### 1.1 Kinds of semantic information

In the following, we briefly describe the kind of semantic information that we investigate in this paper.

**Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

■ Detected sequence: 7, 5, 6, 9      #con = 4, #dis = 2       $\tau = \frac{1}{3}$

■ We obtain

$$\tau_n = \frac{\tau + 1}{2} = \frac{1/3 + 1}{2} = \frac{4/3}{2}$$

**A Benchmark and Evaluation for Text Extraction from PDF**

Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

**ABS 4** Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

**1.1 Kinds of semantic information**  
In the following, we briefly describe the kind of semantic information that we investigate in this paper.

**Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

**A Benchmark and Evaluation for Text Extraction from PDF**

Hannah Bast  
University of Freiburg  
79110 Freiburg, Germany  
bast@cs.uni-freiburg.de

Claudius Korzen  
University of Freiburg  
79110 Freiburg, Germany  
korzen@cs.uni-freiburg.de

**ABS 5** Extracting the body text from a PDF document is an important but surprisingly difficult task. The reason is that PDF is a layout-based format which specifies the fonts and positions of the individual characters rather than the semantic units of the text (e.g., words or paragraphs) and their role in the document (e.g., body text or caption). There is an abundance of extraction tools, but their quality and the range of their functionality are hard to determine.

**1.1 Kinds of semantic information**  
In the following, we briefly describe the kind of semantic information that we investigate in this paper.

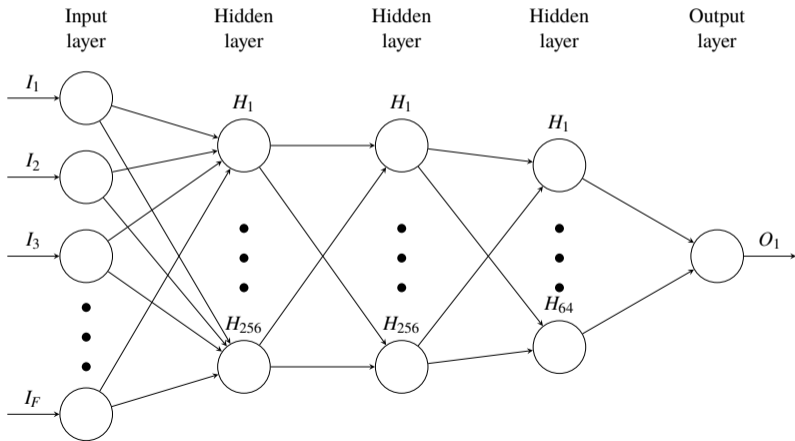
**Word identification.** This is crucial for applications like search: a word that has not been identified correctly will not be found. Word identification in a PDF is non-trivial and challenging for a number of reasons. The spacing between letters can vary from line to line and even within a line, and there is no fixed rule to

■ Detected sequence: 7, 5, 6, 9      #con = 4, #dis = 2       $\tau = \frac{1}{3}$

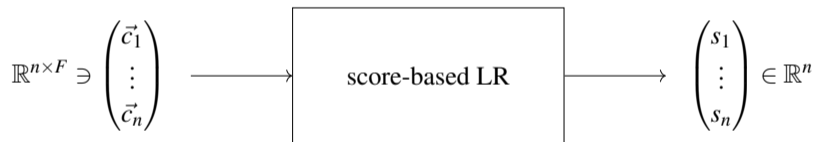
■ We obtain

$$\tau_n = \frac{\tau + 1}{2} = \frac{1/3 + 1}{2} = \frac{4/3}{2} = \frac{2}{3}$$

# Model architecture



- The input is a matrix whose rows correspond to feature representations of cuts



- The output of the model is a score vector
- We then choose the cut with highest score

# Training data format



# name	page num	width,height	subpage	depth	dir
example.pdf	42	612,796	0,0,360,640	2	<b>X</b>
# cut	left/upper semantic roles	right/lower semantic roles	label		
([530,550], <b>Y</b> )	heading	paragraph	1		
([270,290], <b>Y</b> )	paragraph	paragraph	0		
([170,175], <b>X</b> )	-	-	0		
example.pdf	43	612,796	0,0,612,796	1	-
([720,740], <b>Y</b> )	marginal	heading,paragraph	0		
([680,685], <b>Y</b> )	-	-	0		
([420,430], <b>Y</b> )	table,paragraph	formula, caption	0		
([296,316], <b>X</b> )	heading,table,caption	marginal,paragraph,formula	1		



# PdfAct comparison



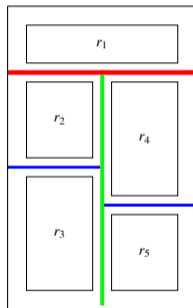
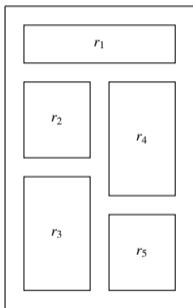
	$B_G^-$	$B_A^-$	$B_G^+$	$B_A^-$	$\tau_n$	$\tau_n^f$
<i>Thesis</i>	51.4%	46.7%	12.9%	14.7%	0.873	0.994
<i>PdfAct</i>	66.5%	54.3%	10.1%	7.5%	0.859	0.985

# Full reading order results



strategy	$\tau_n$	$\tau_n^f$
<i>Largest cut</i>	0.872	0.993
<i>Weighted-largest cut</i>	0.863	0.983
<i>Parameter cut</i>	0.865	0.984
<i>LogisticRegressor</i>	0.873	0.994
<i>BatchClassifier</i>	0.872	0.992
<i>Transformer</i>	0.860	0.978
<i>PdfAct</i>	0.859	0.985

# XY-cut limitations for reading order



- Wang et al.'s LayoutReader shows a way to overcome these limitations