

# VL Search Engines

Lab course: WikiSearch

# Agenda

---

- **Projektbeschreibung**
- Topics
- Prototyp
- Organisatorisches



# Search

HPI

Hasso  
Plattner  
Institut

Page: [1](#) [2](#) [3](#) [4](#) [5](#) [6](#) [7](#) [8](#) [9](#) [10](#)

[Apollo 11](#)

*Snippet for article #662*

[Algorithm](#)

*Snippet for article #775*

[Aircraft](#)

*Snippet for article #849*

[Anna Kournikova](#)

*Snippet for article #890*

[Aluminium](#)

*Snippet for article #900*

search engine

Page: [1](#) [2](#) [3](#) [4](#) [5](#)

[Apollo 11](#)

*Snippet for art*

[Algorithm](#)

*Snippet for art*

[Aircraft](#)

*Snippet for art*

[Anna Kournikova](#)

*Snippet for art*

[Aluminium](#)



Kournikova played at the [Medibank International](#) in Sydney in the first round. She then reached the third round of the [Australian Open](#) (6–4, 4–6, 6–4). She lost in the second round of the [Paris Open](#).

Kournikova reached the semifinals in [Hannover](#). She lost to Novotná in the quarterfinals in [Linz](#), and to [Conchita Martínez](#). She reached her first [WTA Tour](#) final in [Miami](#), where she lost to [Lindsay Davenport](#) (won the first set (2–6, 6–4, 6–1)).

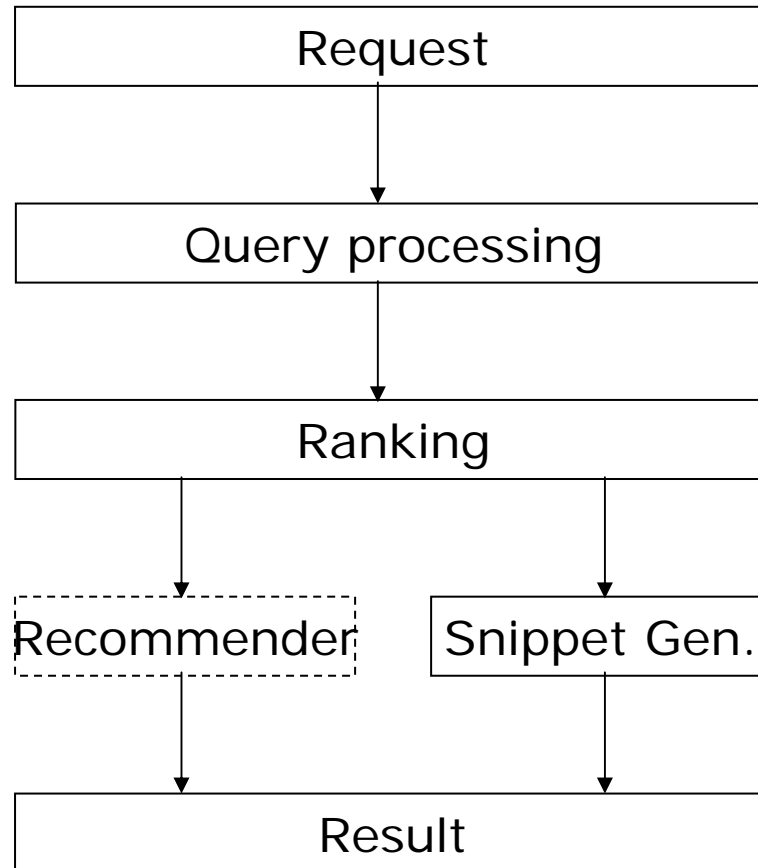
She then played at [Amelia Island](#), when she reached the quarterfinals but lost in the quarter finals of the [Italian Open](#) to [Martina Hingis](#) when she lost to [Conchita Martínez](#). During this tournament she lost to [Hingis](#).

Kournikova then played at the [French Open](#). She lost to [Lindsay Davenport](#) in the semifinals at [Eastbourne](#) (lost to [Sánchez Vicario](#)). She then lost in the first, second, third and fourth rounds. Her last tournament was the [WTA Championships](#), where she lost in the first round to [Monica Seles](#).

In 1998, Kournikova reached her first doubles final, partnering with [Larisa Neiland](#). That was against [Nicola Pietrangeli](#), [Appelmans](#) and [Miriam Oremans](#) in three sets 1–6, 6–3, 7–6(3). They also lost to [Nathalie Tauziat](#). Partnering [Monica Seles](#), she won the Tokyo title. They defeated [Mary Joe Fernandez](#) and [Lindsay Davenport](#). Partnering [Sánchez Vicario](#), she lost to [Lindsay Davenport](#) and [Natasha Zvereva](#) in the finals at [Filderstadt](#) in doubles.<sup>[5]</sup>

**1999**

At the end of the season, Anna Kournikova was ranked #12 in singles and #1 in doubles.<sup>[5]</sup> She was named the world's top female athlete in the world on [Yahoo!](#), the premier search engine of the day.<sup>[6]</sup>



# WikiSearch

- CREATE TABLE **articles** (
  - id int(11) NOT NULL,
  - title text NOT NULL,
  - text text NOT NULL,
  - PRIMARY KEY (id)
 )
  - 8,227,433 Artikel



id	title	text
887	MessagePad	[[Image:Apple Newton MP100.jpg thumb 250px The App...
888	A. E. van Vogt	&lt;!-- Unsourced image removed: [[Image:A.E._van_...
890	Anna Kournikova	{{Infobox tennis player  playername= Anna Kourniko...
891	Accounting	#REDIRECT [[Accountancy]]{{R from related word}}
892	Alfons Maria Jakob	{{Infobox Scientist  name = {{PAGENAME...

# WikiSearch

- CREATE TABLE **inverted\_index** (
  - word varchar(100) NOT NULL,
  - articles text NOT NULL,
  - PRIMARY KEY (word)
 )
  - 7,841,856 Terme



word	
search	303,305,307,308,332,339,569,573,576,584,594,621,628,662,664,681,689,700,736,737,752,775,777,782,791,800,808,824,849,876,890,898
engine	657,662,663,698,775,780,803,848,849,887,890,898,904,974,1009,1014,1022,1027,1132,1146,1175,1202,1208,1271,1311,1336,1358,136

# Team



Prof. Dr. Felix Naumann



Mohammed



Alexander



Jens



Christoph



# Agenda

---

- Projektbeschreibung
- **Topics**
- Prototyp
- Organisatorisches

# Themen

---

- **Core Topics**
  - Index building
  - Query processing
  - Snippet Generation

# Index Construction

---

- **What is an index?**
  - Map from terms to the parts of a document where they occur
- **Why is it important?**
  - Term, document frequency -> weighting
- **Important aspects**
  - Memory, Storage, Processor resources
  - Dynamic indexing
  - Positional indexes
- **The task**
  - Build an index for the Wikipedia collection
- **References**
  - Manning et al., Introduction to Information Retrieval, Chapter 4.
  - Croft et al., Search Engines, Chapter 5

# Index Construction

---

- Relational-Database index.
- Frequencies of each term in each article (term-article basis).
- **Fist Task:**
  - Read related material.
  - Model a relational schema for required tables.
  - Design required classes / modules to scan through all articles and create the index.

# Query processing

- Processing keyword queries
  - Stemming keywords
  - Using an index
  - Query expansion
  
- Result page with relevant articles
  - Simple ranking
  - Paging
  - Spell checking & suggestions
  
- **First Task**
  - Read Croft et al., Search Engines, Section 5.7 & Chapter 6

# Snippet Generation

- **What are Snippets?**
  - A short summary of the document, which is designed to allow the user to decide document's relevance.
- **Its importance?**
  - Help the user determine which result item to choose.
- **Important Aspects:**
  - Static vs. Dynamic snippet generation
- **The task?**
  - Build a snippet generation engine
- **References**
  - Manning et al., Introduction to Information Retrieval, Section 8.7
  - Croft et al., Search Engines, Section 6.2.4

# Snippets Generation

---

- Query-biased snippet generation.
- Caching mechanism?
- **First Task:**
  - Read related material.
  - Evaluate document summarization techniques.

# Themen

---

- **Advanced Topics**
  - Index compression
  - Smart Extract
  - Information extraction
  - Ranking by Paths through the Web
  - Recommendation
  - Search Engine Optimization
  - Search Engine Evaluation
  - Wikify! Learning to Link with Wikipedia



# Index Compression

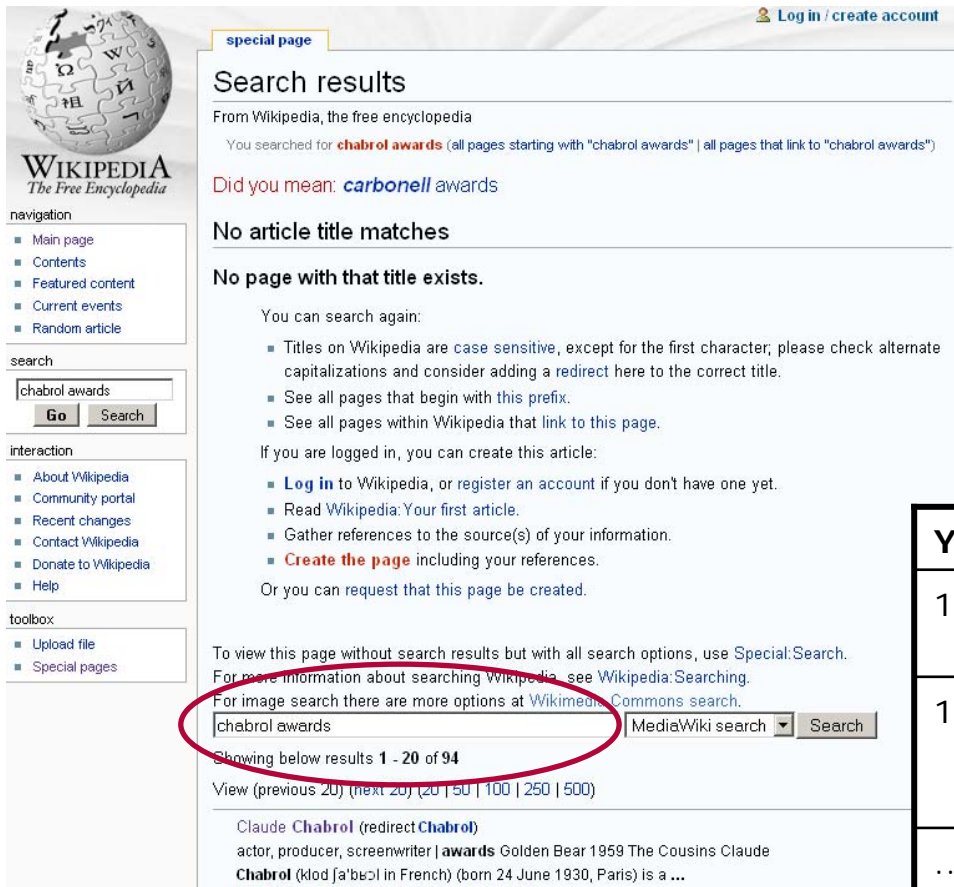
- **Why is compression important?**
  - Less disk space
  - Increased use of caching (frequent terms)
  - Faster transfer of data from disk to memory
- **Important aspects**
  - Fast compression and decompression algorithms
  - Lossless vs. Lossy compression
- **The task**
  - Compress the built index for the Wikipedia collection
- **References**
  - Manning et al., Introduction to Information Retrieval, Chapter 4.
  - Croft et al., Search Engines, Chapter 5

# Index Compression

- Decrease the size of the index by:
  - Stop-words removal.
  - Removal of rare terms, i.e.  $\text{occurrence}(\text{term}) < \text{threshold}$ .
  - Removing redundant data.
  
- **First Task:**
  - Read related material.
  - Compare database- vs. file- based index compression alternatives.
  - Evaluate encoding schemes and choose one.

# Smart Extract (1/3) - Task

- **Extract structured data** from Wikipedia articles (tables) and **find dependencies between them**
- Improved keyword search by using extracted structured data



special page

Search results

From Wikipedia, the free encyclopedia

You searched for **chabrol awards** (all pages starting with "chabrol awards" | all pages that link to "chabrol awards")

Did you mean: [carbonell awards](#)

No article title matches

No page with that title exists.

You can search again:

- Titles on Wikipedia are **case sensitive**, except for the first character; please check alternate capitalizations and consider adding a **redirect** here to the correct title.
- See all pages that begin with **this prefix**.
- See all pages within Wikipedia that **link to this page**.

If you are logged in, you can create this article:

- **Log in** to Wikipedia, or **register an account** if you don't have one yet.
- Read Wikipedia:Your first article.
- Gather references to the source(s) of your information.
- **Create the page** including your references.

Or you can request that this page be created.

To view this page without search results but with all search options, use Special:Search.  
For more information about searching Wikipedia, see Wikipedia:Searching.  
For image search there are more options at Wikimedia Commons search.

chabrol awards  MediaWiki search

Showing below results **1 - 20** of 94

View (previous 20) (next 20) (20 | 50 | 100 | 250 | 500)

Claude **Chabrol** (redirect **Chabrol**)  
actor, producer, screenwriter | **awards** Golden Bear 1959 The Cousins Claude **Chabrol** (klod [a'ʃa'ʁɔl] in French) (born 24 June 1930, Paris) is a ...  
7 KB (981 words) - 22:49, 14 January 2009



article discussion edit this page history

Claude Chabrol

From Wikipedia, the free encyclopedia

**Claude Chabrol** (pronounced [klod [a'ʃa'ʁɔl] in French) (born 24 June 1930, Paris) is a French film director and one of the core members of the **French New Wave** group of filmmakers who first came to prominence in the late 1950s and

Claude Chabrol	
<b>Born</b>	24 June 1930 (age 78) Paris, France
<b>Occupation</b>	director, actor, producer, screenwriter
<b>Years active</b>	1956 - present
Awards won <span>[hide]</span>	
Other awards	
<b>Golden Bear</b>	1959 <i>The Cousins</i>

YEAR	AWARD	MOVIE	...
1959	<a href="#">Berlin International Film Festival</a>	<a href="#">The Cousins (Les cousins)</a>	...
1989	<a href="#">New York Film Critics Circle Award for Best Foreign Language Film</a>	<a href="#">Story of Women</a>	...
...	...	...	...

## Smart Extract (2/3)

- Example: Keyword query *chabrol awards* results in [Claude Chabrol](#), [César Awards 1996](#), [Madame Bovary \(1991 film\)](#), ...
  - Incomplete award list in article of Claude Chabrol (only Golden Bear, 1959)
  - Goal: Offer relevant tables, i.e., awarded movies of Claude Chabrol, next to the well-known result page with the relevant Wikipedia articles
- Wikipedia documents contain lots of structured data
  - Sample Wikipedia articles, containing corresponding relational tables about films: [New York Film Critics Circle Award for Best Foreign Language Film](#), [French films of 1959](#), [San Sebastián International Film Festival](#), [List of crime films: 1980s](#), [List of films set in Las Vegas](#), ...

## Smart Extract (3/3)

---

- **First Task**
  - Read *Uncovering the Relational Web* (WebDB 2008)
  - Start implementing a parser for extracting tables from wikipedia articles

## Information Extraction – IE (1/3)

- Traditional approach
  - Regular expression grammars (e-mails, phone numbers, ...)
  - Drawbacks: Scaling to large data sets and large numbers of rules
- New approach
  - IBM System T
  - An Algebraic Approach to Rule-Based Information Extraction (ICDE 2008)
- **Task 1:** Choose three concepts to extract (companies, persons, ...) from Wikipedia articles and write [AQL](#) scripts for each of these concepts. Evaluate your solution.
- **Task 2:** Extend the keyword search with a semantic autocompletion feature/advanced query language using your extracted concepts.

# Information Extraction – IE (2/3)

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Proin elementum non ante. **John Pipe** played the **guitar**. Aliquam erat volutpat. Curabitur a massa. ~~Vivamus luctus, risus in sagittis facilis, arcu augue rutrum ve~~

Regex match

0-30 characters

Dictionary match

<b>&lt;ProperNoun&gt;</b>	<b>&lt;within 30 characters&gt;</b>	<b>&lt;Instrument&gt;</b>
Regular Expression Match		Dictionary Match

```

select
  CombineSpans(name.match, instrument.match)
  as annot
from
  Regex(/[A-Z]\w+(\s[A-Z]\w+)?/, DocScan.text) name,
  Dictionary("instr.dict", DocScan.text) instrument
where
  Follows(0, 30, name.match, instrument.match);
  
```

(c) Declarative Information Extraction, The Avatar Group IBM Almaden Research Center, 2008

# Information Extraction – IE (3/3)

- **First Task**

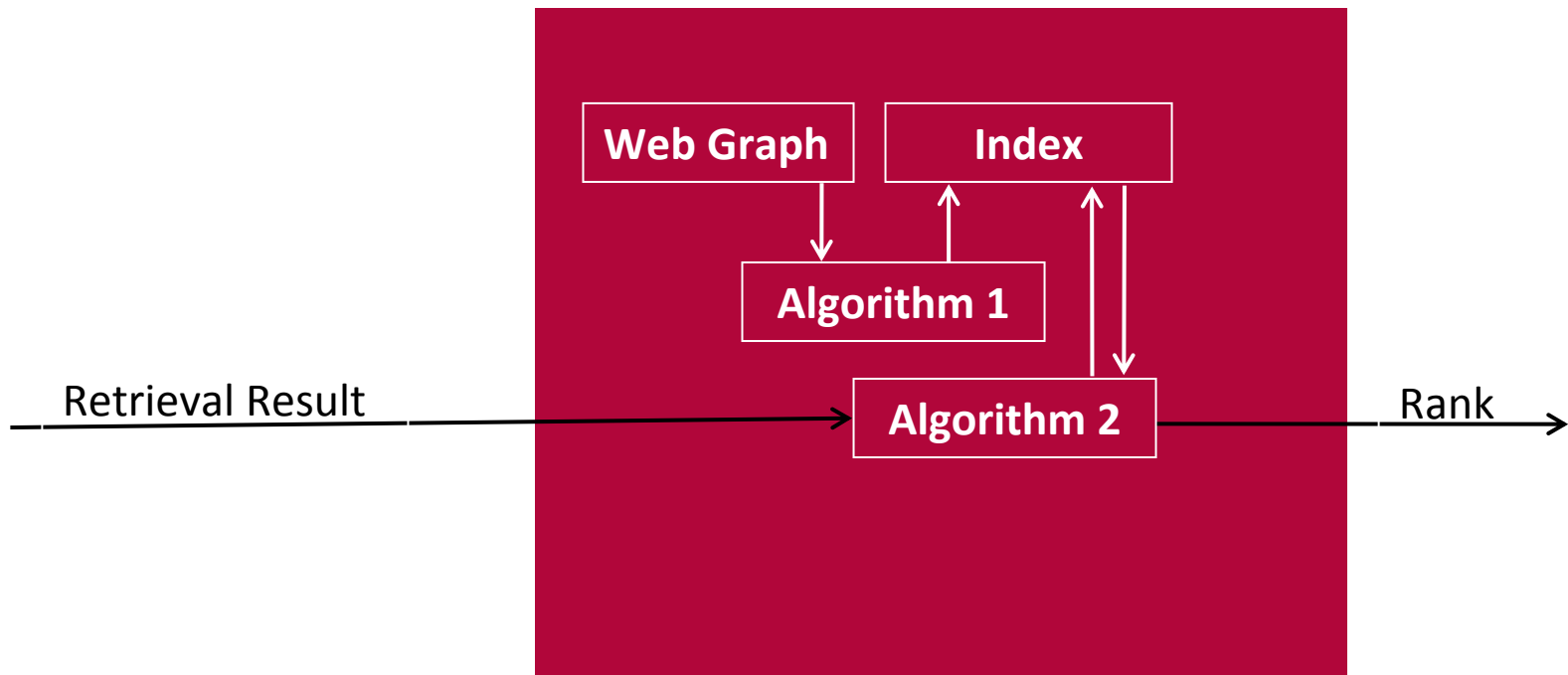
- Read *An Algebraic Approach to Rule-Based Information Extraction* (ICDE 2008)
- Install IBM System T
- Start writing [AQL](#) scripts for wikipedia articles



## Ranking by Paths through the Web (1/2)

- SPRINT: Search by Path Ranks on the INTranet
- Idea:
  - Imagine you were a computer scientist
  - You search for ‚genetic‘ on Wikipedia
  - ‚genetic programming‘ should be ranked higher than ‚genetics‘ or ‚genetic diversity‘
  - Unless you are a biologist
- Implementation: Rank retrieval results according to the (link) distance from a user defined central webpage
- Prerequisite:
  - We know who you are or a webpage that fits your characteristics best
- **First task:** Read *Efficient search ranking in social networks*, Vieira et al., CIKIM 2007

## Ranking by Paths through the Web (2/2)



# Recommendation

---

- Expand the results of a search by attaching recommendations
  - Find one or more “recommended item” to every search
  - Items are other Wikipedia Pages
  - Recommendation should be based on interactions with pages, or other data that is not already used in a ranking
- Two possible ways of recommending items:
  - Use authorship information: authors that have written parts of this page have also written these other pages
  - Use usage/view information: users that have seen (clicked on) this results have also seen these other pages

## Recommendation – first steps

- Get familiar with recommendation systems:
  - Read two easy articles (get a print copy from Jens, next week)
- Think about page interactions, and answer the questions:
  - What interactions are possible? What data for which interactions is available? What data do you want to use?
  - How do you want to gather and store the data? Think about a schema, plan gathering the data.
- Collect first small amounts of interaction data (a sample):
  - Download data
  - Crawl data, screen scrape data
  - Create or collect data by hand
- Meet me at May, 14<sup>th</sup>, with that data and the plan

# Search Engine Evaluation

- Which search engine answers Wikipedia keyword queries best?
  - Wikipedia
  - Google
  - Yahoo
  - ...
- Compare the quality of results and the ranking over a defined set of queries
  - Precision & recall
- Identify and evaluate applied techniques, such as stemming, query expansion, ...
- Repeat your experiments within certain periods (weeks) and determine possible changes
- **First task:** Read *Croft et al., Search Engines, Chapter 8*

# Wikify! Learning to Link with Wikipedia

- Automatically recognize topics mentioned in unstructured text and link them to the appropriate Wikipedia articles
- Machine learning approach
  - Can be used to identify significant terms within unstructured text
  - Millions of manually-defined links found within Wikipedia articles are used as training set
- **First task:** Read *Learning to Link with Wikipedia*, CIKM 2008 & *Wikify! Linking Documents to Encyclopedic Knowledge*, CIKM 2007

**Iranian POW negotiator holds talks with Iraqi ministers**

The head of [Iran's prisoner of war](#) commission met with two [Iraqi](#) Cabinet ministers Saturday in a bid to glean information about thousands of Iranian POWs allegedly in Iraq, the official Iraqi News Agency reported.

Iraqi Foreign Minister [Mohammed Saeed al-Sahhaf](#) told Abdullah al-Najafi that the two states needed to "speed up the closure of what remains from the POW and Missing-in-Action file," INA said.

The issue of POWs and missing persons remains a stumbling block to normalizing relations between the two neighbors.

Iraq has long maintained that it has released all Iranian prisoners captured in the [1980-88 Iran-Iraq War](#). The countries accuse each other of hiding POWs and preventing visits by the [International Committee of the Red Cross](#) to prisoner camps.

The ICRC representative in [Baghdad](#), Manuel Bessler, told [The Associated Press](#) that his organization has had difficulty visiting POWs on both sides on a regular basis.

In April, Iran released 5,584 since [1990](#).

More than 1 million people w

**Baghdad**

Baghdad is the capital of Iraq and of Baghdad Governorate. With a metropolitan area estimated at a population of 7,000,000, it is the largest city in Iraq. It is the second-largest city in the Arab world (after Cairo) and the second-largest city in southwest Asia (after Tehran).

[open in wikipedia](#)

# Agenda

---

- Projektbeschreibung
- Topics
- **Prototyp**
- Organisatorisches

# Prototyp

---

R:\SearchEngines09\_Naumann\Materialien\_nosync\

- MySQL installieren
- Wikipedia-Datenbank anlegen und laden (siehe HowTo)
- Apache Tomcat installieren & anpassen (siehe HowTo)
- WikiSearch Anwendung entpacken und starten



# Agenda

---

- Projektbeschreibung
- Topics
- Prototyp
- **Organisatorisches**

# Organisatorisches

---

- Liste mit 3 Themenwünschen per Email an Alexander Albrecht (mindestens eine Core Topic).
- E-Mail bis zum 26.04.2009, 23:59 Uhr
- Vergabe der Themen und Zuordnung der Teams am 27.04.
- Teams vereinbaren individuelles Treffen mit ihren Betreuern bis spätestens 14.05. (Terminvorschläge)
- Bewertung der implementierten Lösung