



Hasso  
Plattner  
Institut

IT Systems Engineering | Universität Potsdam



## Search Engines – Exercise

### Smart Data Extraction

Franz Goerke

Thomas Berger

# Our guideline

2

- Requirements:
  - Extraction of tables / structured data
  - Analyze table structures
  - Enable search for tables
  - Tableindex creation
  
- Enhancements:
  - Get relation among tables
  - Link resolving
  - Statistics about table contents

# New Database structure

3

## ■ New Tables

- WikiTables
- ArticleTitles

TableID	ArticleID	Content
ArticleID	Title	

## ■ New Indices

- Inverted Table Index
- Linked Table Index
- Attribute Table Index

TableID	ArticleID	Content
TableID	LinkedTables	BacklinkedTables
Attribute	Tables	

# Relations among tables

4

- Approaches
  - Link analysis
    - Tables contain links to articles
      - » Not very useful
      - » Links mostly refer to just one single term
  - Content analysis
    - For every term query inverted index
    - Count occurrences of tables for every word
    - Use top results
    - Very useful
  - Attribute analysis
    - Not useful
    - Most attribute names are equal (name, year, language)

# Linked tables

Oscar Mathisen Award winners

Year	Winner	Nationality
1959	<a href="#">Knut Johannesen</a>	<a href="#">Norway</a>
1960	<a href="#">Boris Stenin</a>	<a href="#">Soviet Union</a>
1961	<a href="#">Henk van der Grift</a>	<a href="#">Netherlands</a>
1962	<a href="#">Jonny Nilsson</a>	<a href="#">Sweden</a>
1963	<a href="#">Nils Aaen</a>	<a href="#">Norway</a>
1964	<a href="#">Ants Antson</a>	<a href="#">Soviet Union</a>
1965	<a href="#">Per Ivar Moe</a>	<a href="#">Norway</a>

1. Extract Links

via Bliki API

article\_titles

id	title
10	<a href="#">AccessibleComputing</a>
15	<a href="#">AfghanistanPeople</a>
20	<a href="#">AfghanistanMilitary</a>
25	<a href="#">Autism</a>
30	<a href="#">AsWeMayThink</a>
35	<a href="#">AlbaniaGovernment</a>

lookup

tables

tableid	articleid	content
1	39	{  class=&quot;wikitable&quot; border=&quot;1&quot;...
2	290	{  class=&quot;wikitable&quot; - style=&quot;back...
3	290	{  cellspacing=&quot;10&quot; cellpadding=&quot;0&quot;...
4	303	{  class=&quot;wikitable&quot; border=&quot;1&quot;...
5	303	{  class=&quot;wikitable&quot; border=&quot;1&quot;...
6	305	{  class=&quot;wikitable&quot; border=&quot;1&quot;...
7	324	{  class=&quot;sortable wikitable&quot; style=&quot;...

Opt: Backlinking

tableid	linkedTables	backLinkedTables
4	332059,17500	NULL
7	14064,218503,218504,218505,218506	NULL
10	101533,101534	NULL
12	4354,4355,4356,4464,4465,4466,4467	NULL

Thanks for your Attention

6

# Demo & Questions