



Hasso
Plattner
Institut

IT Systems Engineering | Universität Potsdam

Link Analysis goes MapReduce

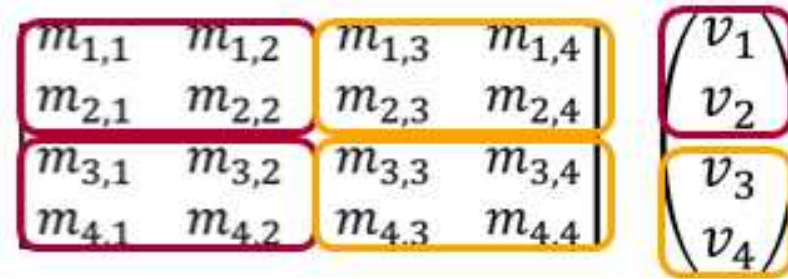
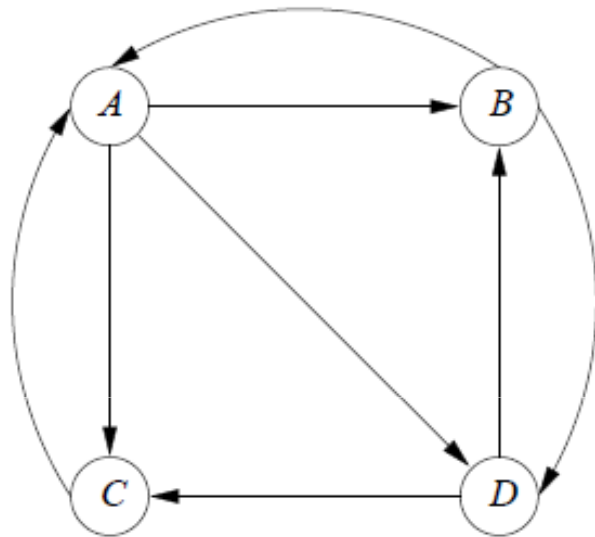
The Hadoop Implementation

Thomas Zimmermann

Philipp Berger

Flashback

2

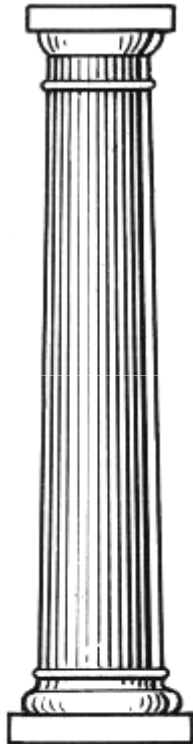


Innsbruck

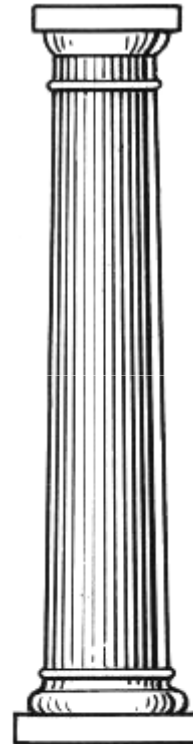
Country	Austria
State	Tyrol
Administrative region	Statutory city
Population	117,342 (2006)
Area	104.91 km ²
Population density	1,119 /km ²
Elevation	574 m
Coordinates	47°16' N 11°23' E
Postal code	6010-6080
Area code	0512
Licence plate code	I
Mayor	Hilde Zach
Website	www.innsbruck.at

Overview

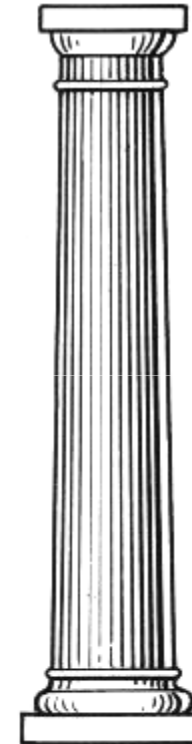
3



1. Pre- / Postprocessing



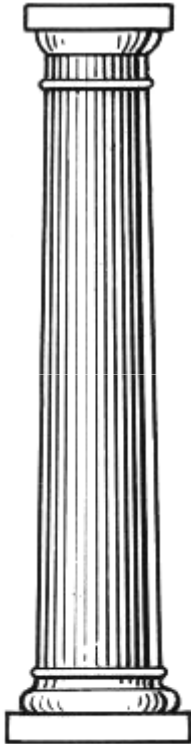
2. Our Jobs



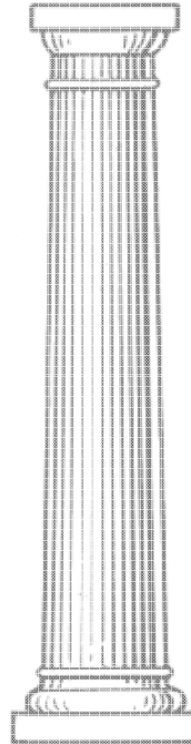
3. Evaluation

Overview

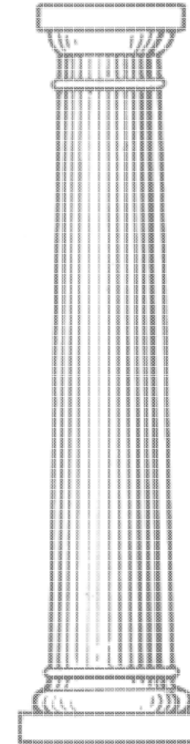
4



1. Pre- / Postprocessing



2. Our Jobs



3. Evaluation

Our Dataset

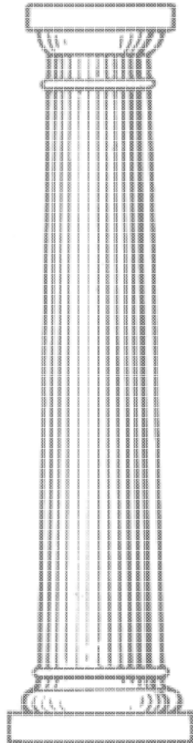
5

- DBPedia Infobox Properties

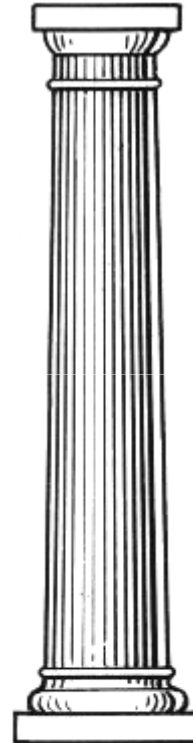


Overview

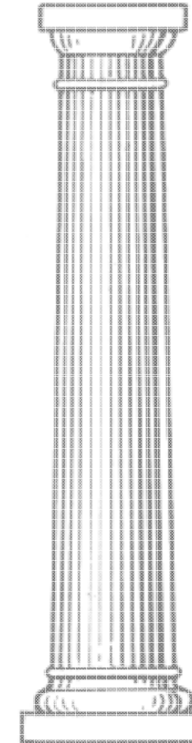
6



1. Pre- / Postprocessing



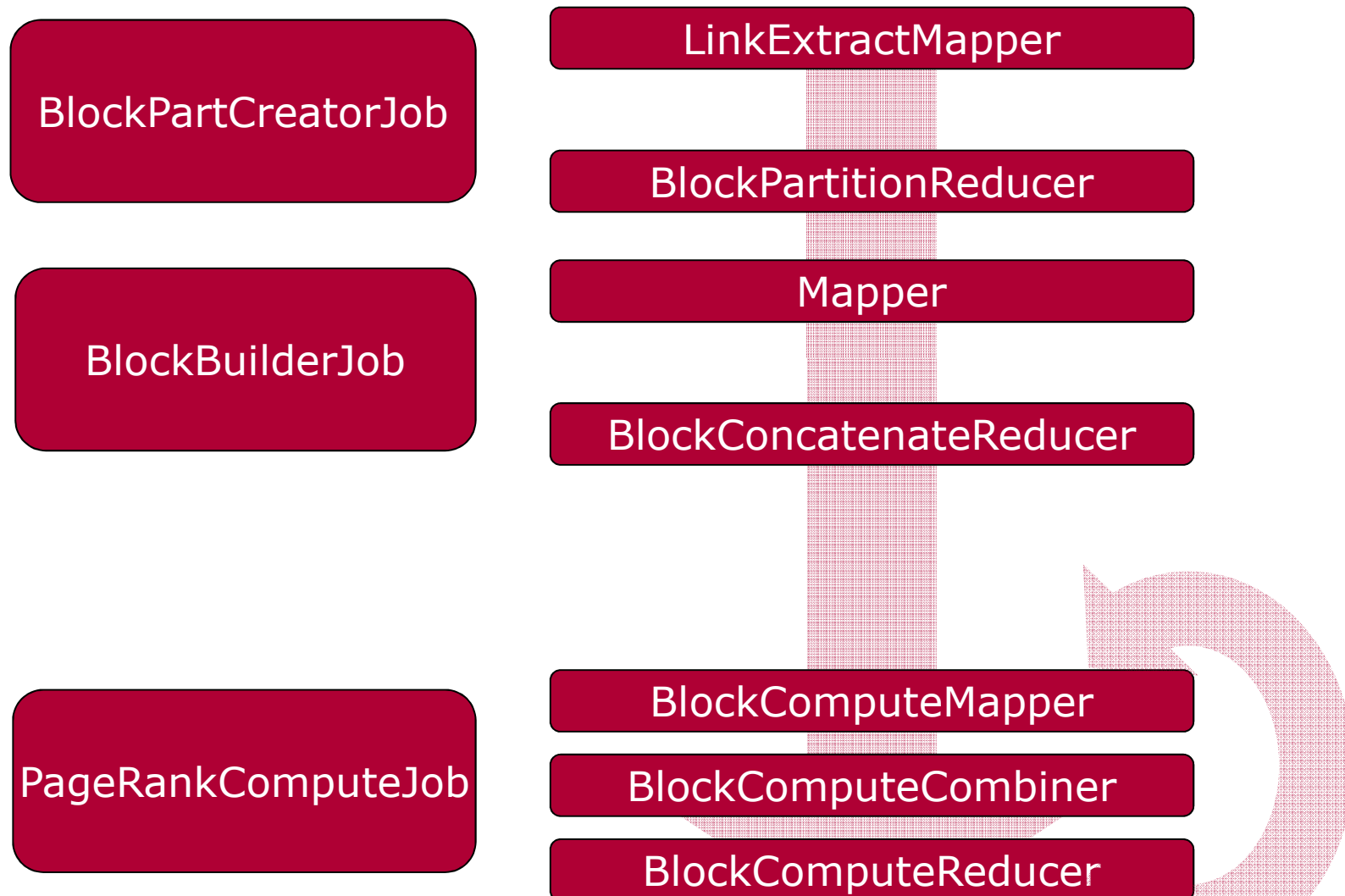
2. Our Jobs



3. Evaluation

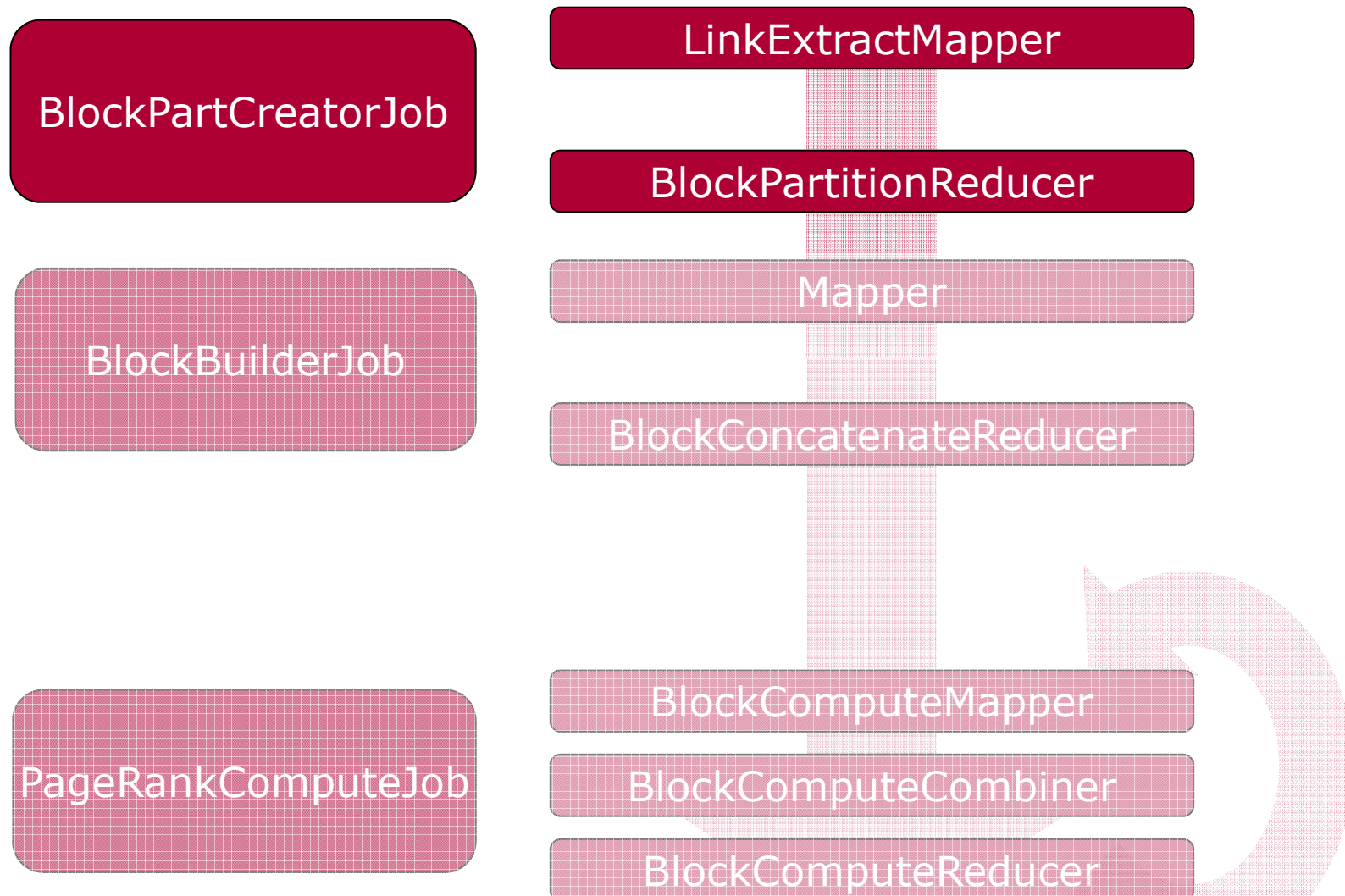
Our Jobs

7



Our Jobs

8



Our Jobs – Create Block Parts (Map)

9

LinkExtractMapper

Input:

Key, (32_(Eiffel Tower), 122_(France))

Output:

32_(Eiffel Tower), 122_(France)

Our Jobs – Create Block Parts (Reduce)

10

BlockPartitionReducer

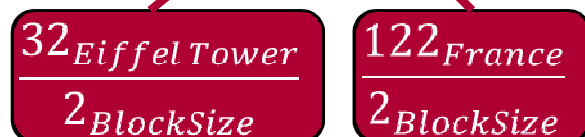
	32 _{Eiffel Tower}	33 _{Sushi}
122 _{France}	X	
123 _{Japan}		X
124 _{Paris}	X	

Input:

$32_{\text{Eiffel Tower}}, [122_{\text{France}}, 124_{\text{Paris}}]$

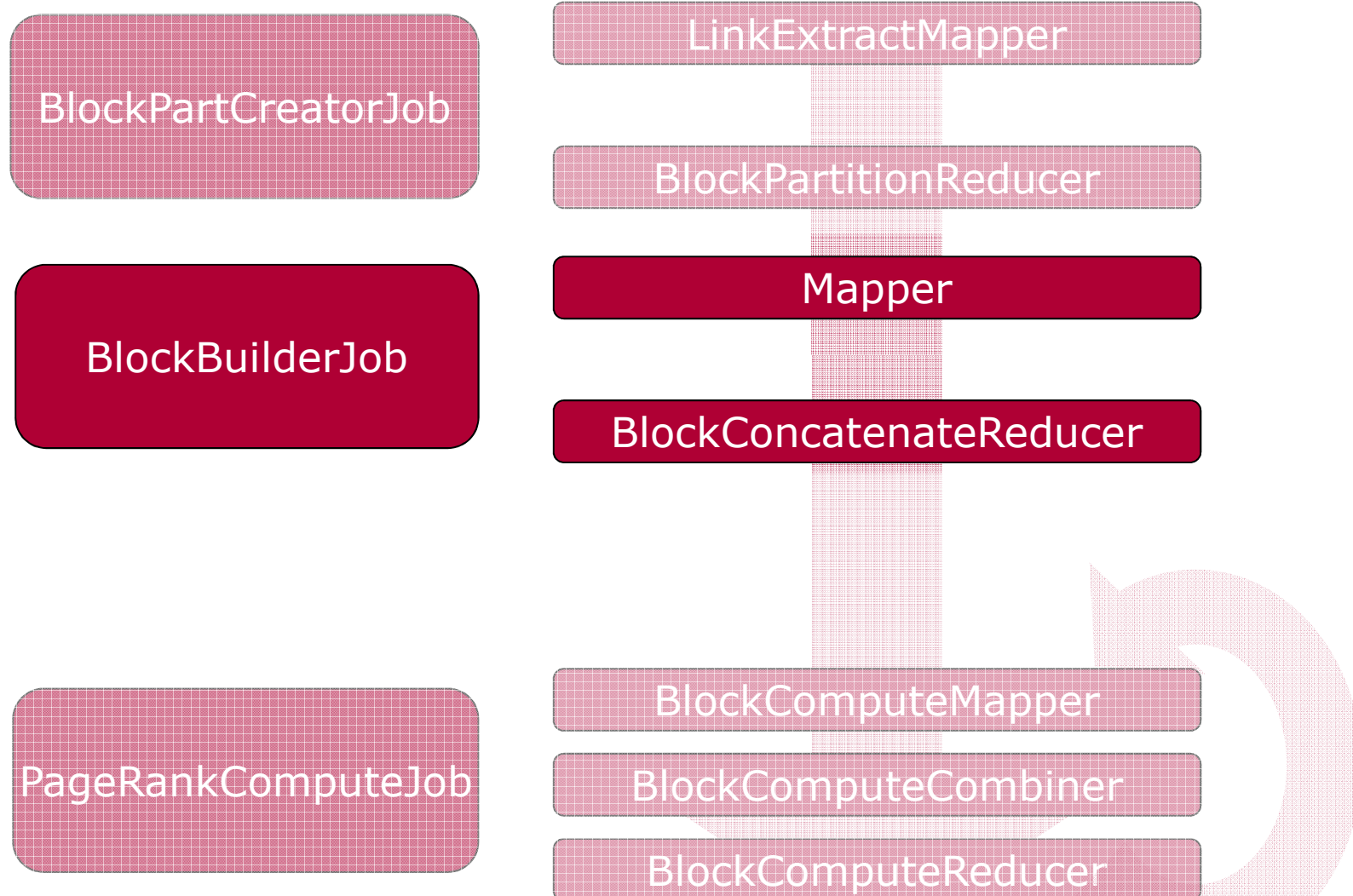
Output:

$(16, 61), (32_{\text{Eiffel Tower}}, 122_{\text{France}}, 2 \text{ #OutgoingLinks})$



Our Jobs

11



Our Jobs – Concatenate Block Parts

12

BlockConcatenateReducer

	32 _{Eiffel Tower}	33 _{Sushi}
122 _{France}	X	
123 _{Japan}		X
124 _{Paris}	X	

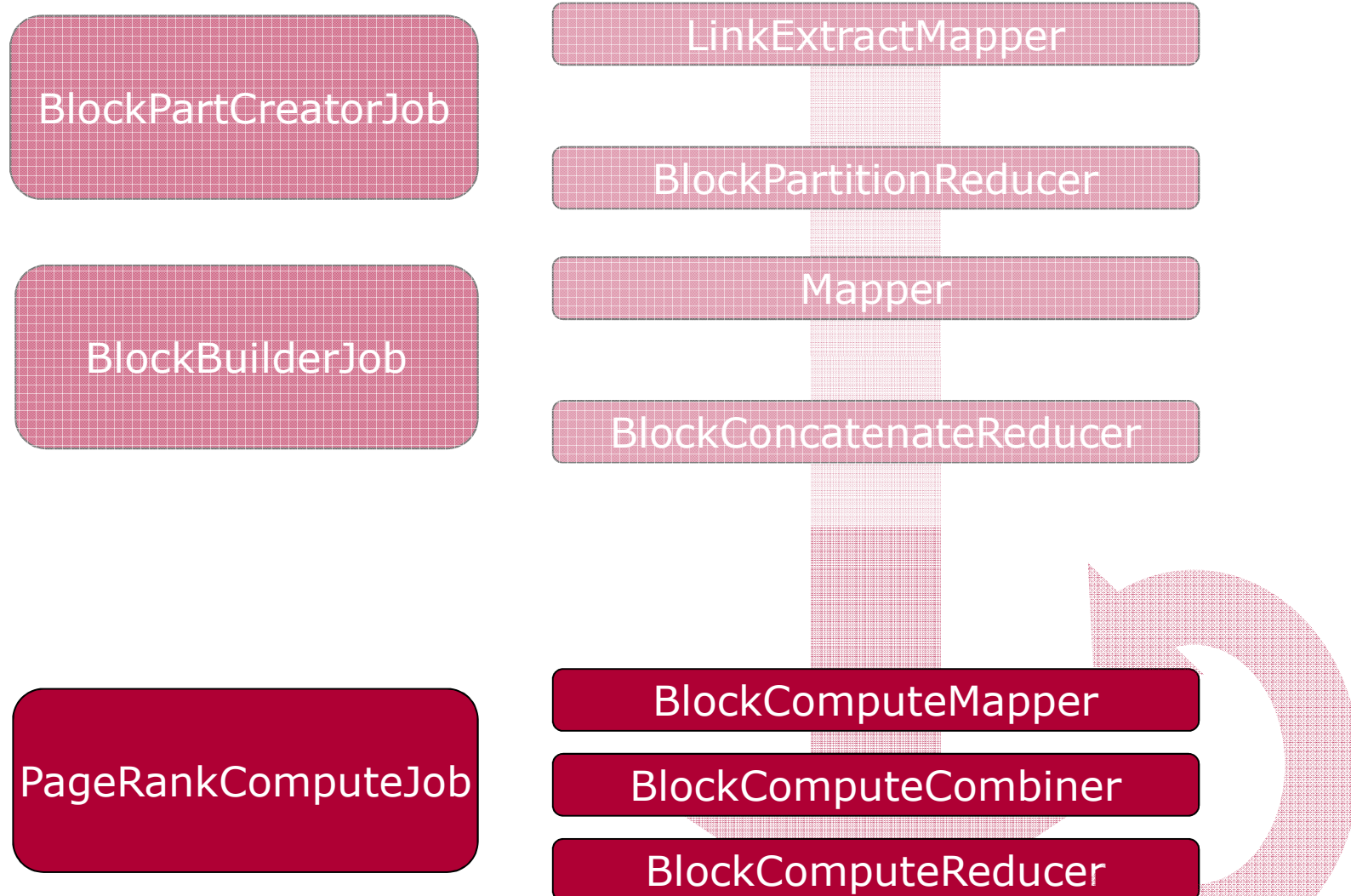
Input / Output:

(16_{Eiffel Tower,Sushi}, 61_{France,Japan}),

[(32_{Eiffel Tower}, 122_{France}, 2_{#OutgoingLinks}), (33_{Sushi}, 123_{Japan}, 1_{#OutgoingLinks})]

Our Jobs

13



Our Jobs – Do PageRank

14

BlockComputeMapper

Input:

$(16_{\text{Eiffel Tower, Sushi}}, 61_{\text{France, Japan}})$,

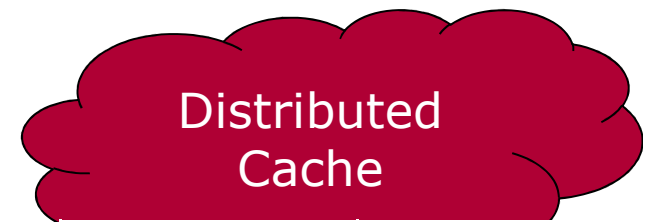
$((32_{\text{Eiffel Tower}}, 122_{\text{France}}, 2_{\text{\#OutgoingLinks}}), (33_{\text{Sushi}}, 123_{\text{Japan}}, 1_{\text{\#OutgoingLinks}}))$

- Get PageRank of sources for one Block

Output:

$122_{\text{France}}, \text{rank}(\text{Eiffel Tower}) * 1/2$

$123_{\text{Japan}}, \text{rank}(\text{Sushi}) * 1/1$



Page	Rank
32 _{Eiffel Tower}	0.125
33 _{Sushi}	0.125
34 _{Island}	0.125

Our Jobs – Do PageRank

15

BlockComputeCombiner

Input:

$122_{\text{France}}, \text{rank}(\text{Eiffel Tower}) * \frac{1}{2}$

$122_{\text{France}}, \text{rank}(\text{Paris}) * \frac{1}{4}$

- Simply sums up incoming values

Output:

$122_{\text{France}}, \text{rank}(\text{Eiffel Tower}) * \frac{1}{2} + \text{rank}(\text{Paris}) * \frac{1}{4} \}$

Our Jobs – Do PageRank

16

BlockComputeReducer

Input:

$122_{\text{France}}, \text{rank}(\text{Eiffel Tower}) * \frac{1}{2} + \text{rank}(\text{Paris}) * \frac{1}{4}$

$122_{\text{France}}, \text{rank}(\text{Lyon}) * \frac{1}{4}$

```
// Random Jump
```

```
result = beta * pageRankSum + (1 - beta) * (1 / pageCount);
```

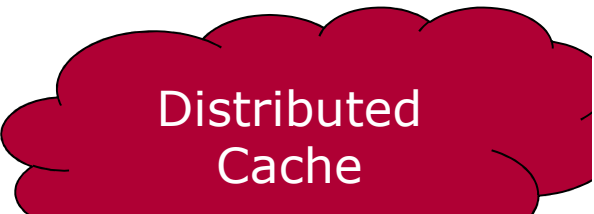
Output:

$122_{\text{France}}, \text{result}$

Our Jobs – Do PageRank

17

- Distributed Cache Update
 - Read the output of PageRank calculation
 - Load the PageRank into the Distributed Cache
- Start next iteration with new PageRank in Distributed Cache
 - Until **convergence** is reached
 - Until the root-mean-square-deviation < target variance

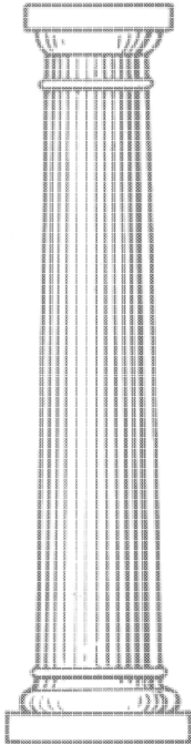


Page	Rank
32 _{Eiffel Tower}	0.105
33 _{Sushi}	0.025
34 _{Island}	0.177

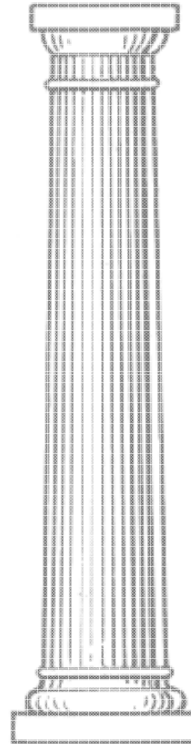


Overview

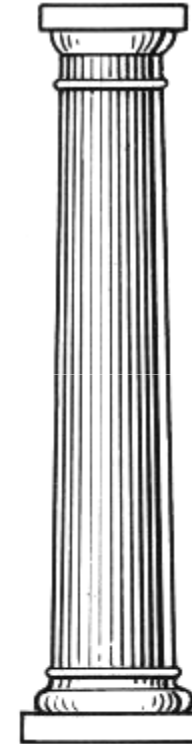
18



1. Pre- / Postprocessing



2. Our Jobs



3. Evaluation

Results: Top 10 German Sites

19

Rank	PageRank	Site
1	0.0084	< http://www.opengis.net/gml/_Feature >
2	2.9127E-4	< http://dbpedia.org/resource/Rheinland-Pfalz >
3	2.4293E-4	< http://dbpedia.org/resource/Deutschland >
4	2.3606E-4	< http://dbpedia.org/resource/Vereinigte_Staaten >
5	1.4299E-4	< http://dbpedia.org/resource/Piemont >
6	1.4109E-4	< http://dbpedia.org/resource/Schleswig-Holstein >
7	1.3978E-4	< http://dbpedia.org/resource/Texas >
8	1.3291E-4	< http://dbpedia.org/resource/United_States >
9	1.3159E-4	< http://dbpedia.org/resource/Baden-Wuerttemberg >
10	1.2901E-4	< http://dbpedia.org/resource/Niedersachsen >

Results: Top 10 German Sites

20

- 2477 pages link to „Rheinland-Pfalz“!

Liste der Städte und Gemeinden in Rheinland-Pfalz

Das [deutsche Bundesland Rheinland-Pfalz](#) besteht aus insgesamt

- 2.306 politisch selbstständigen Städten und Gemeinden (Stand: 7. Juni 2009).

- 1242 to „Piemont“ (Italy’s largest region)

Kategorie:Ort im Piemont

In dieser Kategorie werden alle 1206 Gemeinden der [italienischen Region Piemont](#) geführt.

Rank	Site
1	< http://www.opengis.net/gml/_Feature >
2	< http://dbpedia.org/resource/Rheinland-Pfalz >
3	< http://dbpedia.org/resource/Deutschland >
4	< http://dbpedia.org/resource/Vereinigte_Staaten >
5	< http://dbpedia.org/resource/Piemont >
6	< http://dbpedia.org/resource/Schleswig-Holstein >
7	< http://dbpedia.org/resource/Texas >
8	< http://dbpedia.org/resource/United_States >
9	< http://dbpedia.org/resource/Baden-Württemberg >
10	< http://dbpedia.org/resource/Niedersachsen >

Evaluation

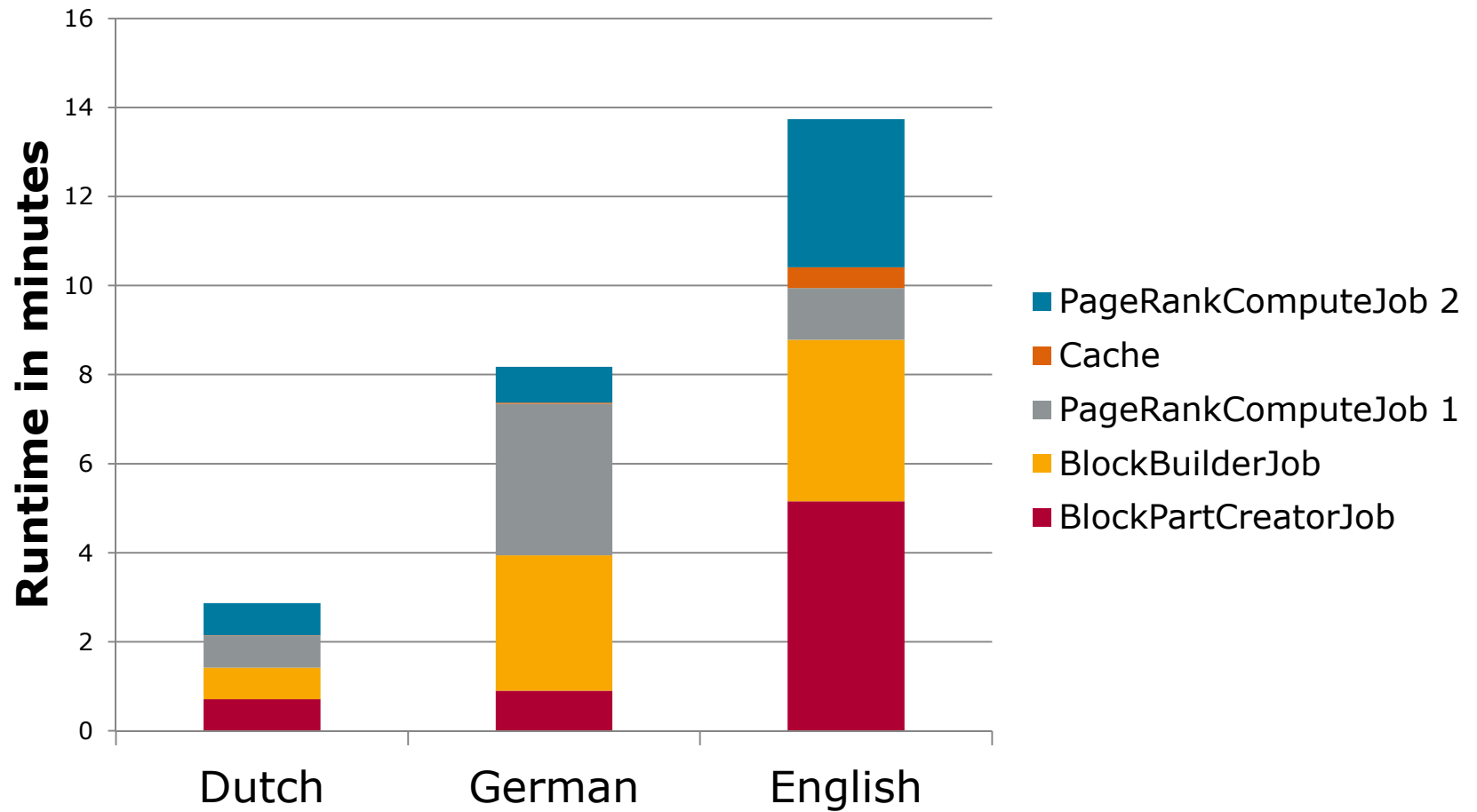
21

Data	SingleThreaded	Cluster
Dutch 1.9 MB	2 sec	2.53min
German 154.8 MB	4 min	8.25min
English 2.2 GB	>3 Std	14.22min

- Single-Threaded graph analysis library
 - `edu.uci.ics.jung.algorithms.scoring.PageRank`
- On Cluster with 10 nodes
 - Blocksize 10 000

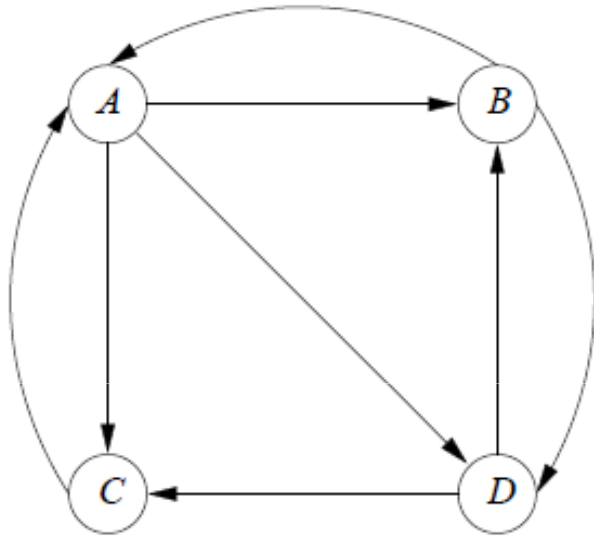
Evaluation

22



Summary

23



$m_{1,1}$	$m_{1,2}$	$m_{1,3}$	$m_{1,4}$	v_1
$m_{2,1}$	$m_{2,2}$	$m_{2,3}$	$m_{2,4}$	
$m_{3,1}$	$m_{3,2}$	$m_{3,3}$	$m_{3,4}$	v_3
$m_{4,1}$	$m_{4,2}$	$m_{4,3}$	$m_{4,4}$	

BlockPartCreatorJob

BlockBuilderJob

PageRankComputeJob