



**Hasso
Plattner
Institut**

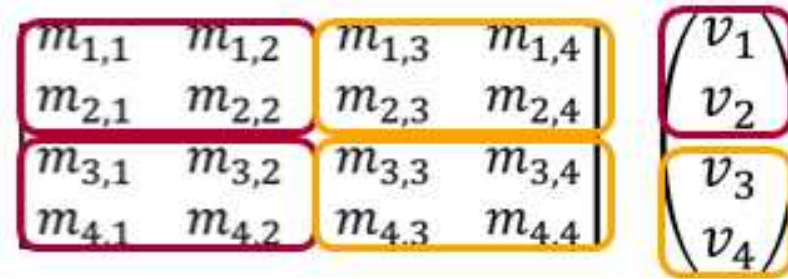
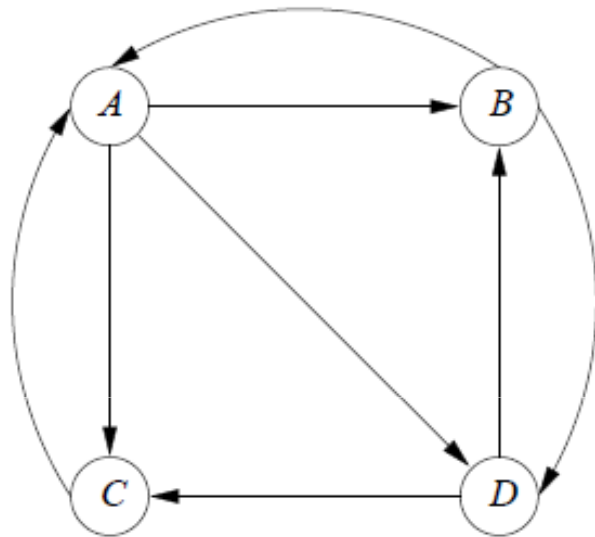
IT Systems Engineering | Universität Potsdam

Link Analysis goes to the Stratosphere

Thomas Zimmermann
Philipp Berger

Recap

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

Agenda

3

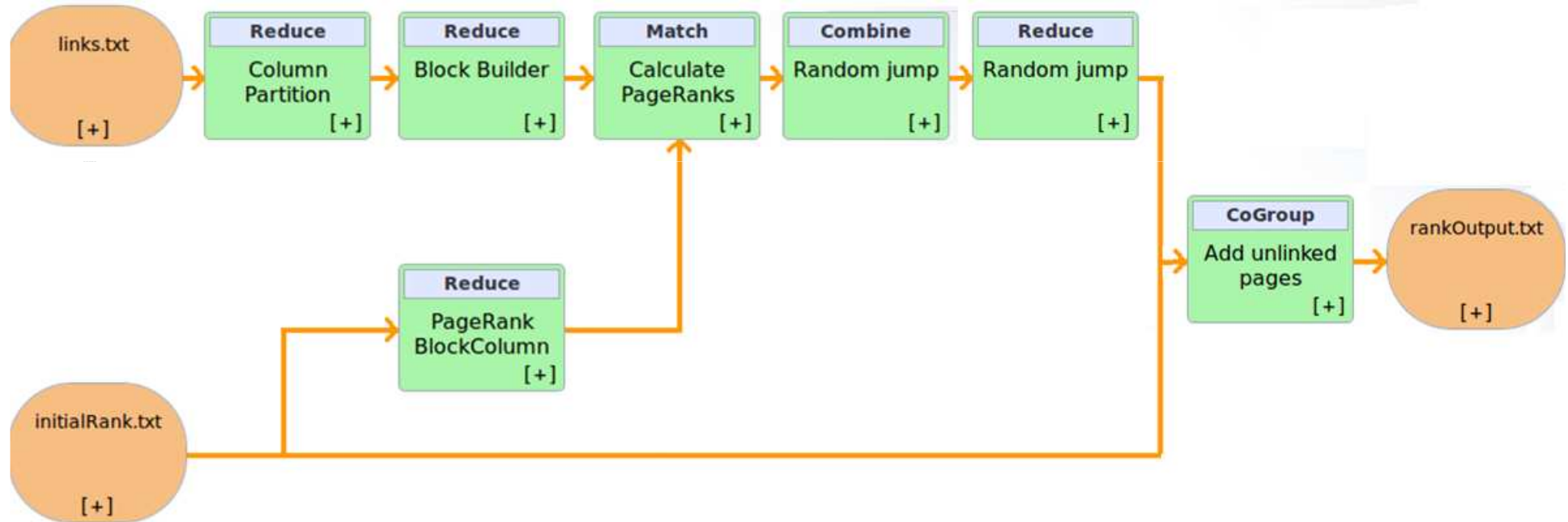
1. PageRank on Stratosphere
2. Hadoop vs. Stratosphere
3. Evaluation

4

- 1. PageRank on Stratosphere**
2. Hadoop vs. Stratosphere
3. Evaluation

PageRankPlan

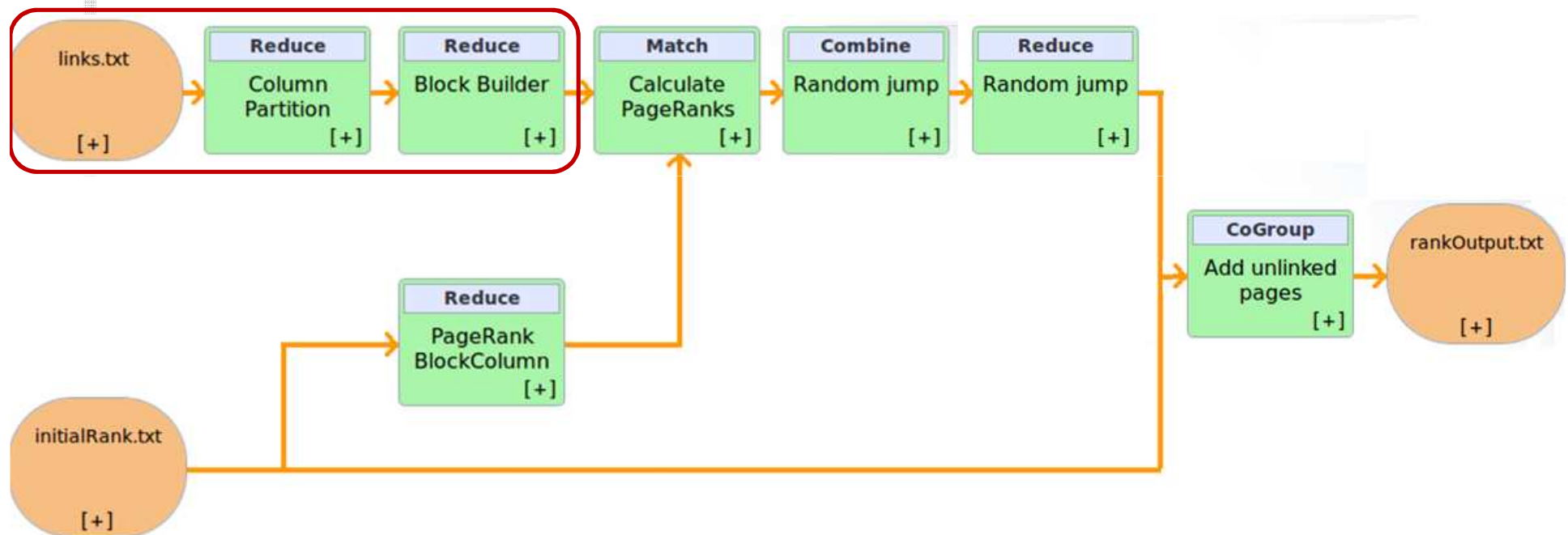
5



PageRankPlan

6

➤ Block Building



InputFormat Links

7

Innsbruck



Country	Austria
State	Tyrol
Administrative region	Statutory city
Population	117,342 (2008)
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

Record

Source Page Id

Target Page Id

BlockBuilding

8

➤ Column Partitioner Reducer

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

Record
Block Column Id
Block Row Id
BlockColumnValue

```
class BlockColumnValue implements Value
```

```

long sourceSiteId;
Collection<Long> targetSiteIds;
long targetSetSize;

```


BlockBuilding

9

➤ Column Partitioner Reducer

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

Record	
Block Column Id	16
Block Row Id	61
BlockColumnValue	(32, [122], 2)

```
class BlockColumnValue implements Value
```

```

long sourceSiteId;
Collection<Long> targetSiteIds;
long targetSetSize;

```

BlockBuilding

10

- BlockPartConcatenate

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

Record
Block Column Id
Block Row Id
BlockColumnValueList

```
class BlockColumnValueList extends PactList<BlockColumnValue>
```

BlockBuilding

11

- BlockPartConcatenate

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

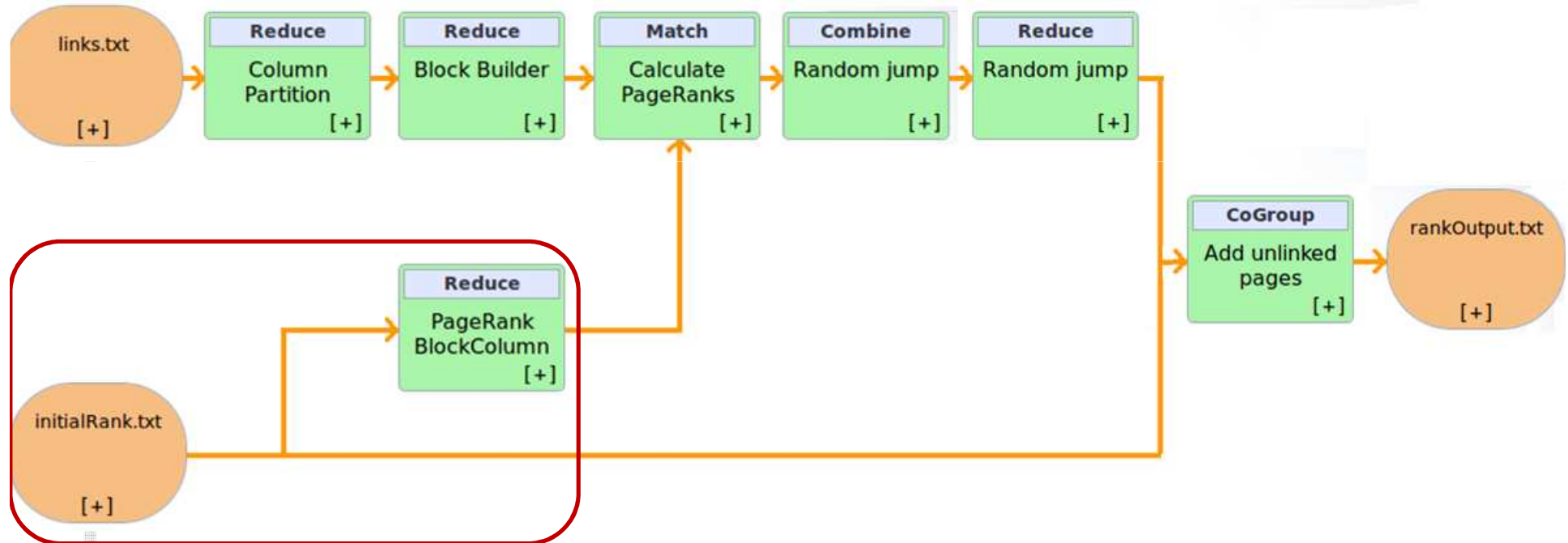
Record	
Block Column Id	16
Block Row Id	61
BlockColumnValueList	(32, [122], 2), (33, [123], 1)

```
class BlockColumnValueList extends PactList<BlockColumnValue>
```

PageRankPlan

12

➤ Read Input Ranks



Inputformat Ranks

13

- Assign block column id to pages
 - To group page ranks for each block

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

Record
Block Column Id
Page Id
Rank

Inputformat Ranks

14

- Assign block column id to pages
 - To group page ranks for each block

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

Record	
Block Column Id	16
Page Id	33
Rank	0.125

RankColumn Reducer

15

- Creates Map for each Block

Record	
Block Column Id	16
Page Id	32
Rank	0.125

Record	
Block Column Id	16
Page Id	33
Rank	0.125

Record	
Block Column Id	
PageToRankMap	

```
class PageToRankMap extends PactMap<PactLong, PactDouble>
```

RankColumn Reducer

16

➤ Creates Map for each Block

Record	
Block Column Id	16
Page Id	32
Rank	0.125

Record	
Block Column Id	16
Page Id	33
Rank	0.125

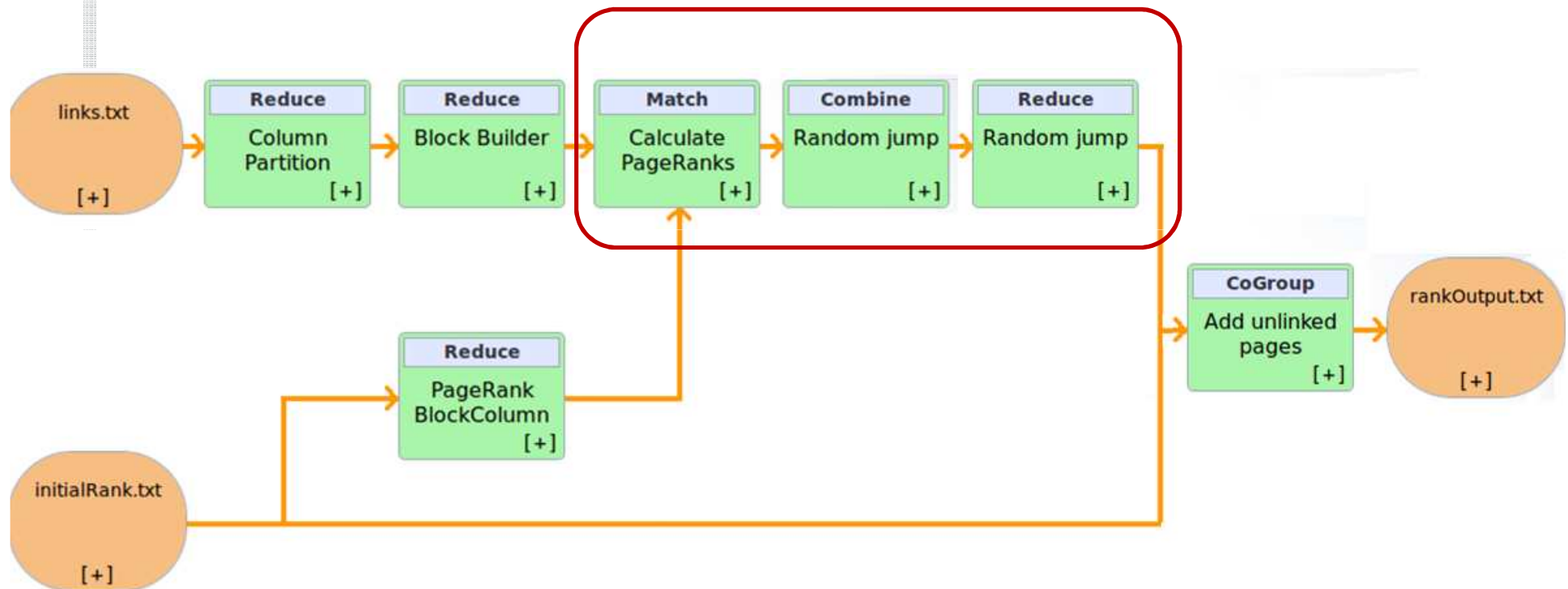
Record	
Block Column Id	16
PageToRankMap	32 → 0.125 33 → 0.125

```
class PageToRankMap extends PactMap<PactLong, PactDouble>
```


PageRankPlan

17

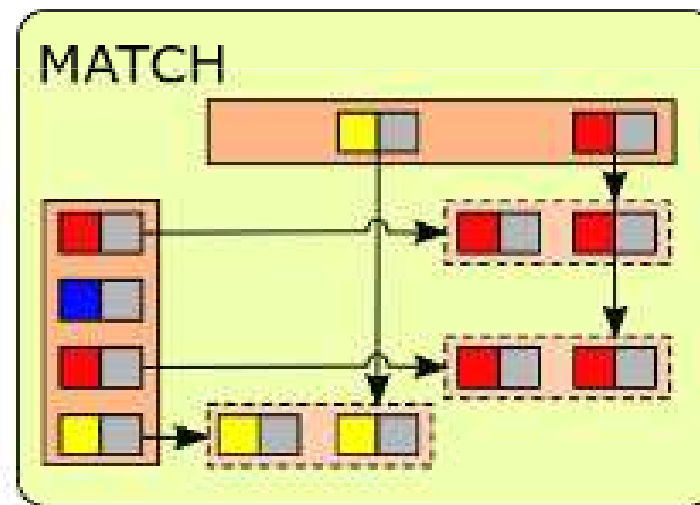
➤ Calculate PageRank



Stratosphere Match

18

- Like a database JOIN

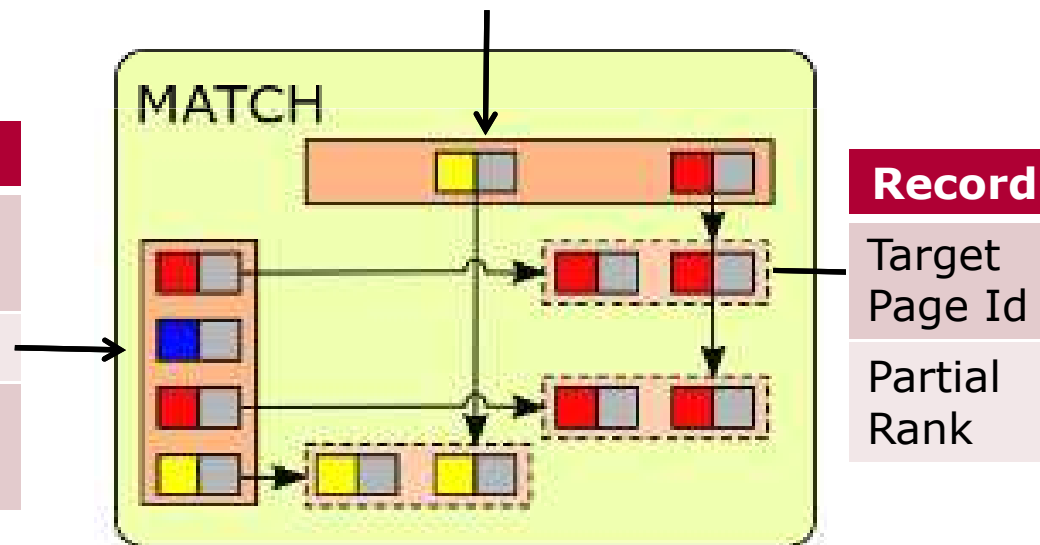


Match Rank and Block

19

PageToRankMap	
Block Column Id	16
PageToRankMap	32 → 0.125 33 → 0.125

Record	
Block Column Id	16
Block Row Id	61
BlockColumn ValueList	(32,[122],2), (33,[123],1)

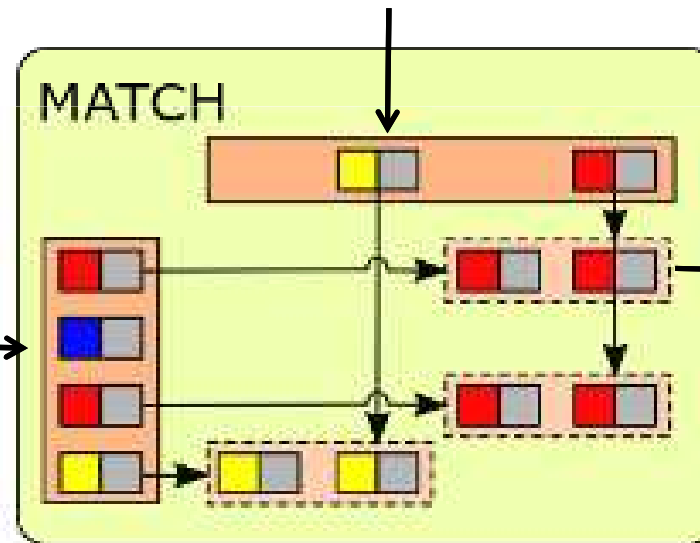


Match Rank and Block

20

PageToRankMap	
Block Column Id	16
PageToRankMap	32 → 0.125 33 → 0.125

Record	
Block Column Id	16
Block Row Id	61
BlockColumn ValueList	(32,[122],2), (33,[123],1)



Record	
Target Page Id	122
Partial Rank	0.0625

Reduce to get final Rank

21

- RandomJump
 - Combine
 - Sums Up new ranks
 - Reduce
 - Calculates the final sum
 - Adds the RandomJump component

Record

Target Page Id

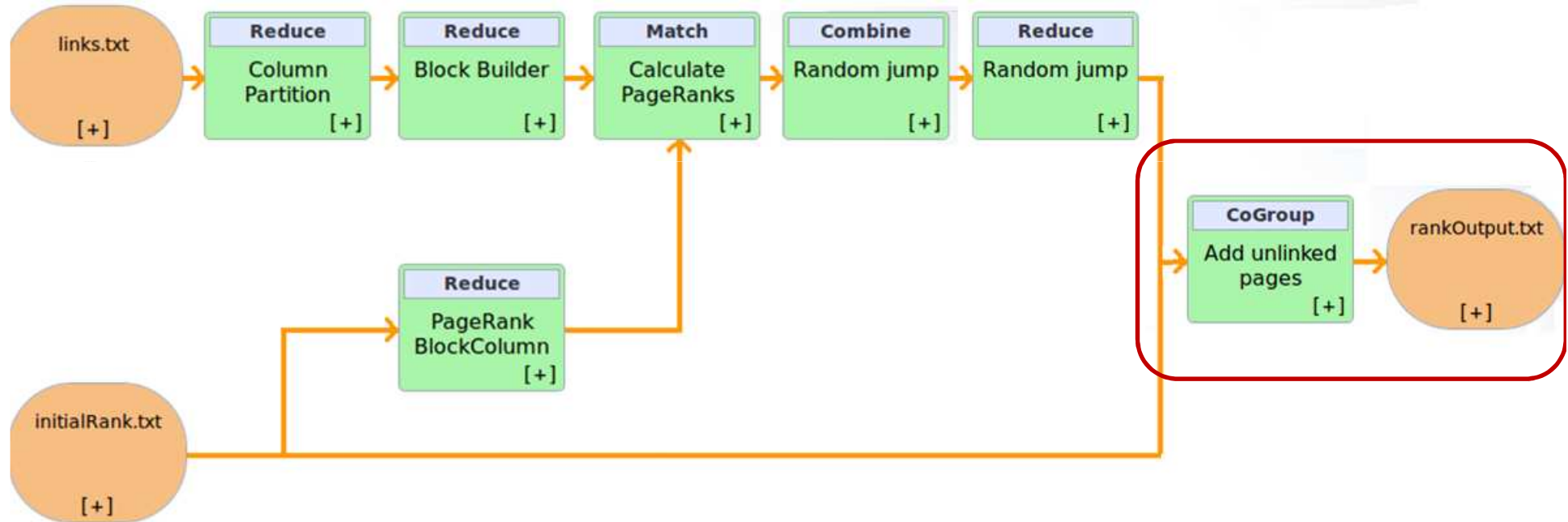
Calculated Rank

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

PageRankPlan

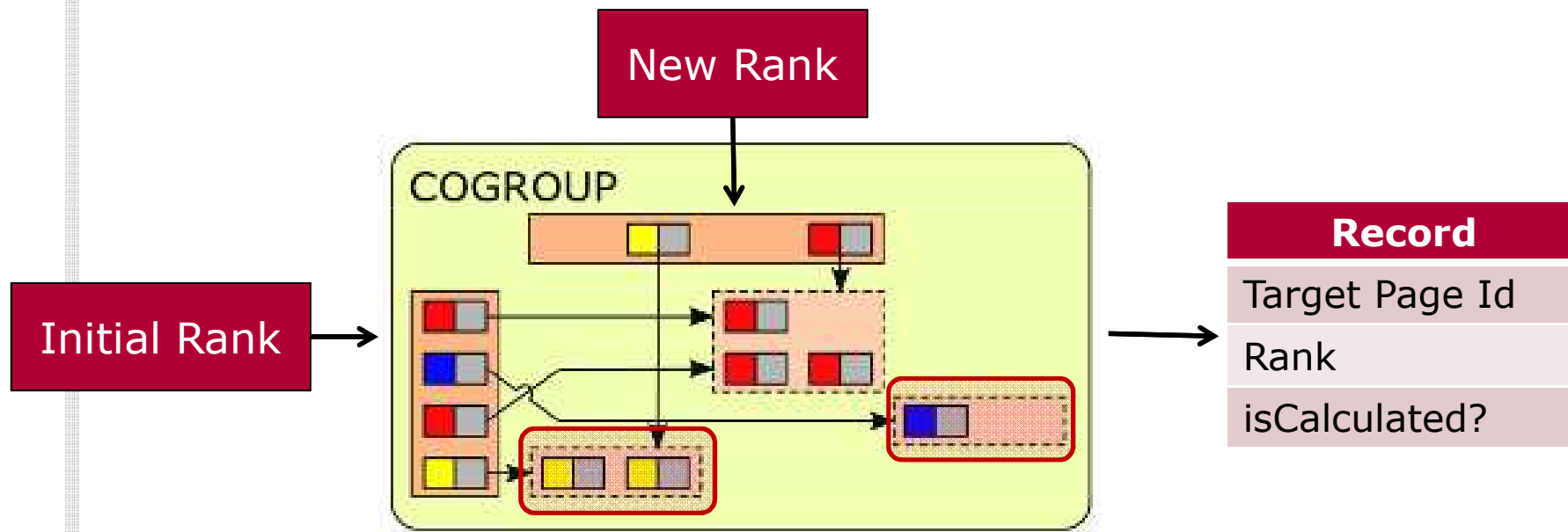
22

➤ Calculate PageRank



Some ranks are missing?

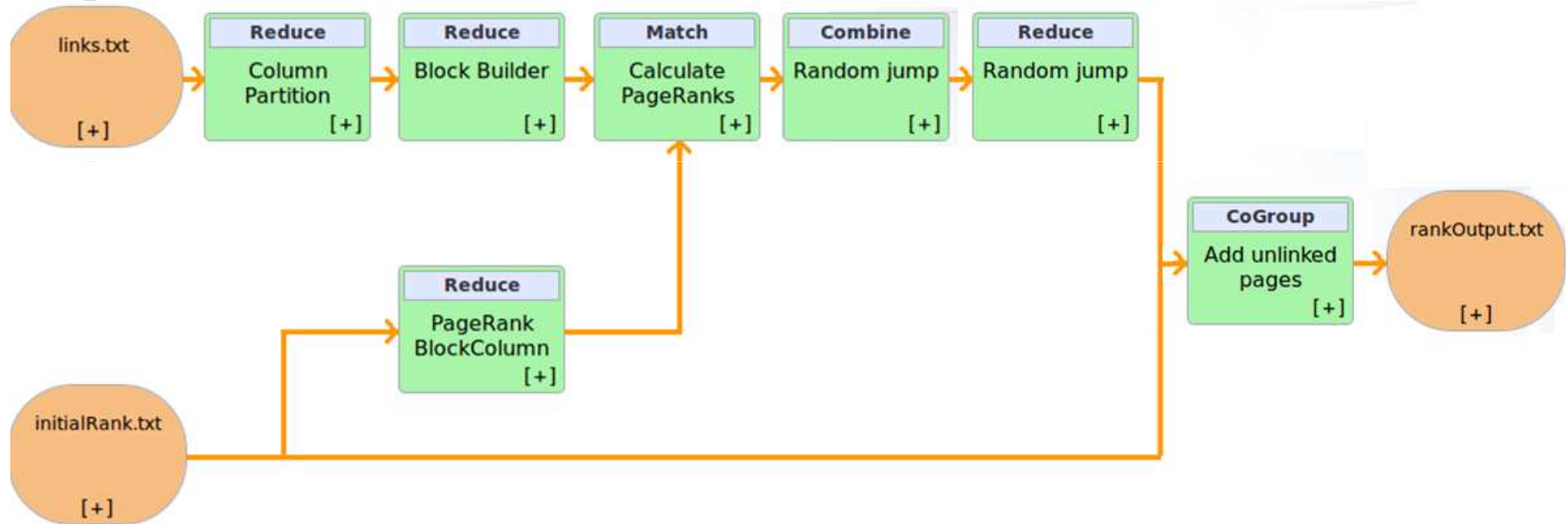
23



➤ $(InitialRanks \setminus NewRanks) \cup NewRanks$

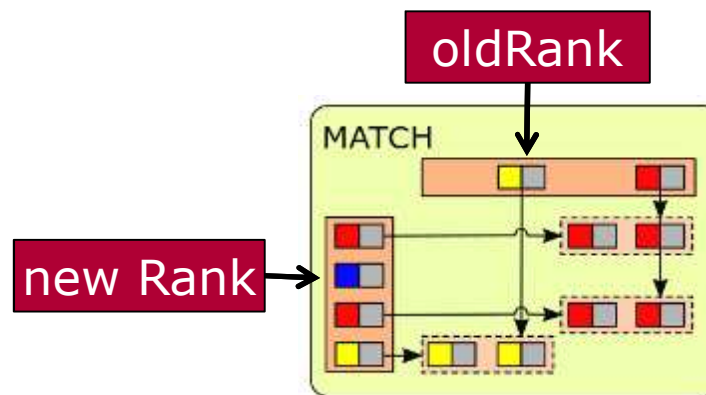
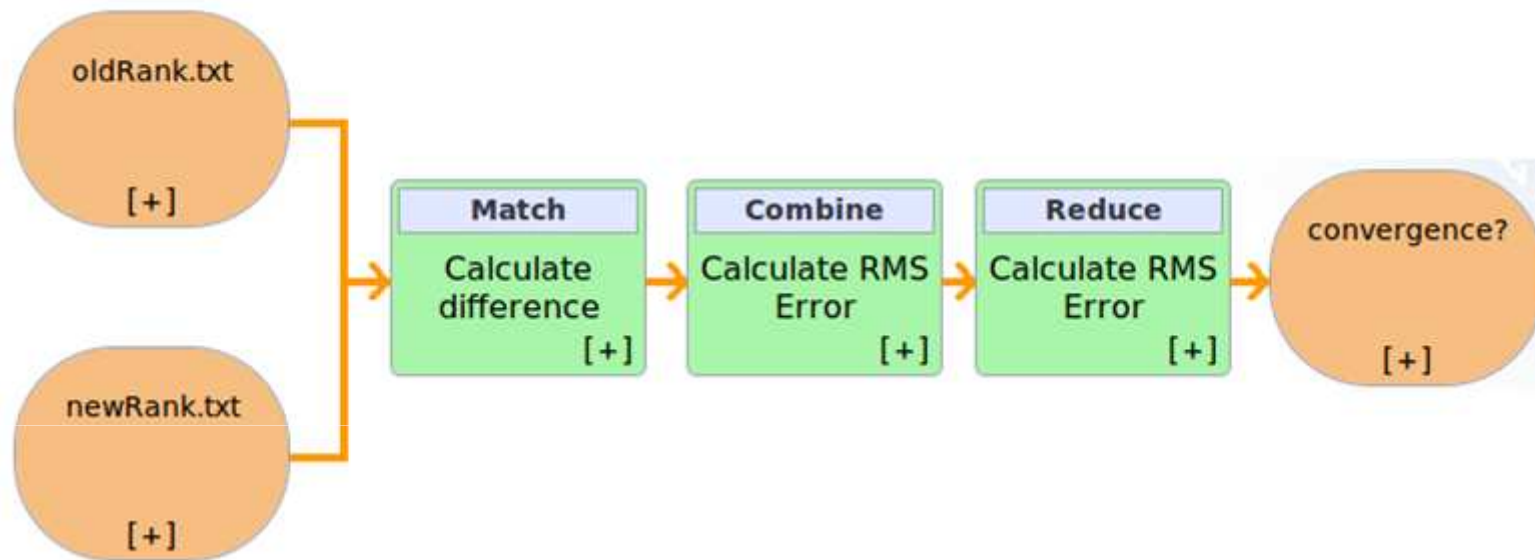
PageRankPlan

24



Does it converge?

25



Result Ranking

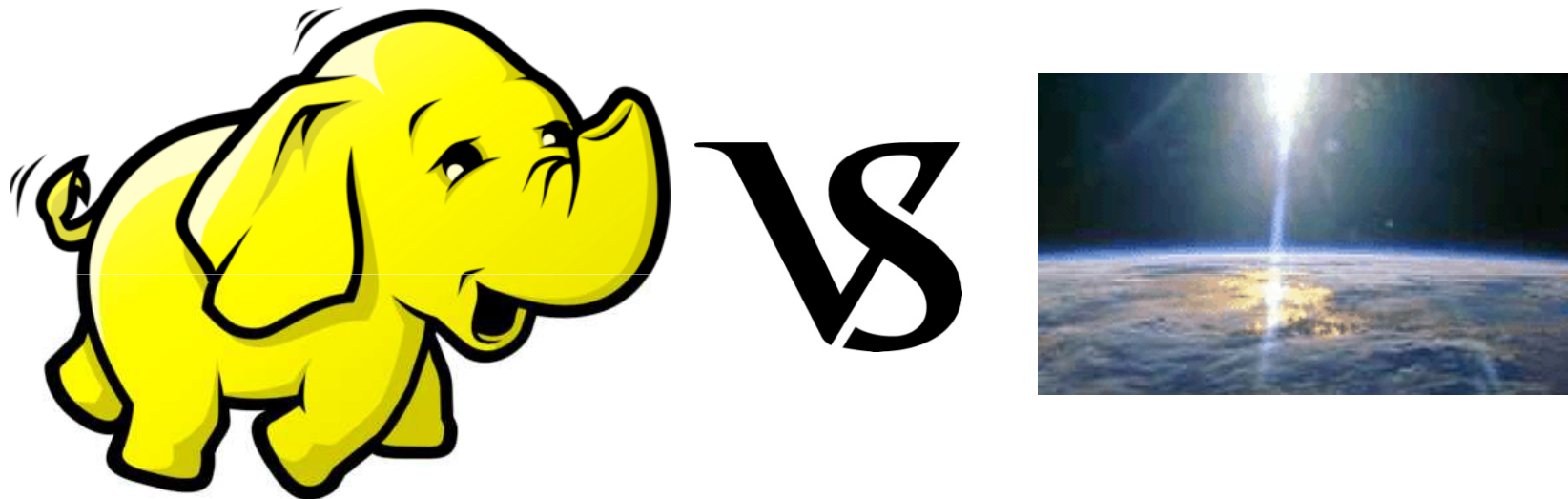
26

Rank	PageRank	Site
1	0.00972794215888403	< http://www.opengis.net/gml/ _Feature >
2	0.004640235231407156	< http://dbpedia.org/resource/United_States >
3	0.0033554866451856414	< http://dbpedia.org/resource/Eukaryote >
4	0.0019618452796001246	< http://dbpedia.org/resource/Animal >
5	0.0013024558125835565	< http://dbpedia.org/resource/France >
6	0.0011861730946363563	< http://dbpedia.org/resource/United_Kingdom >
7	0.001178379282062479	< http://dbpedia.org/resource/English_language >
8	0.001100648265149723	< http://dbpedia.org/resource/Plant >
9	0.0010822226849594	< http://dbpedia.org/resource/Poland >
10	0.0008418068734001761	< http://dbpedia.org/resource/Australia >

1. PageRank on Stratosphere
- 2. Hadoop vs. Stratosphere**
3. Evaluation

Hadoop vs. Stratosphere

28



➤ *from the coders point of view*

InputFormat and Mapper

29

Hadoop

- Need for an extra mapper for input file splitting
- Complex structure of InputFormat in combination with RecordReader and InputSplit

Stratosphere

- Easy adaptable InputFormat interface
- Definition of keys in the contract

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	6010-6060
Area code	0512
Licence plate code	I
Mayor	Hilde Zach
Website	www.innsbruck.at

Type Safety

30

- Hadoop
 - Types of key and value at compile time via Generics
 - TupleWritable
- Stratosphere
 - Types of record fields only known at runtime
 - Combined keys



Distributed Cache and Match

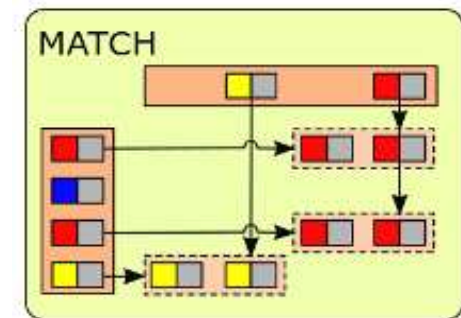
31

Hadoop

- Distributed Cache spreads source ranks to each mapper
- Each mapper reads from same Cache

Stratosphere

- Combine multiple inputs with Match
- Match gets the subset of ranks for one block



Need for Initial Rank

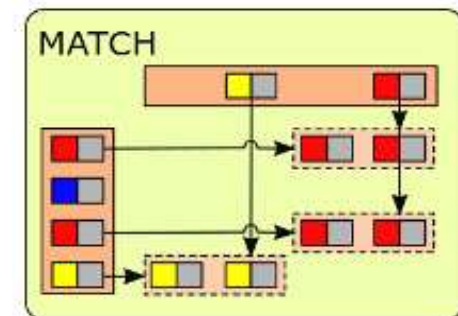
32

Stratosphere with Match

- explicit initial rank
- Co-Group for unlinked pages

Hadoop with Distributed Cache

- On-the-fly default page rank calculation



Stability

33

- Hadoop
 - Windows support
 - Jobtracker

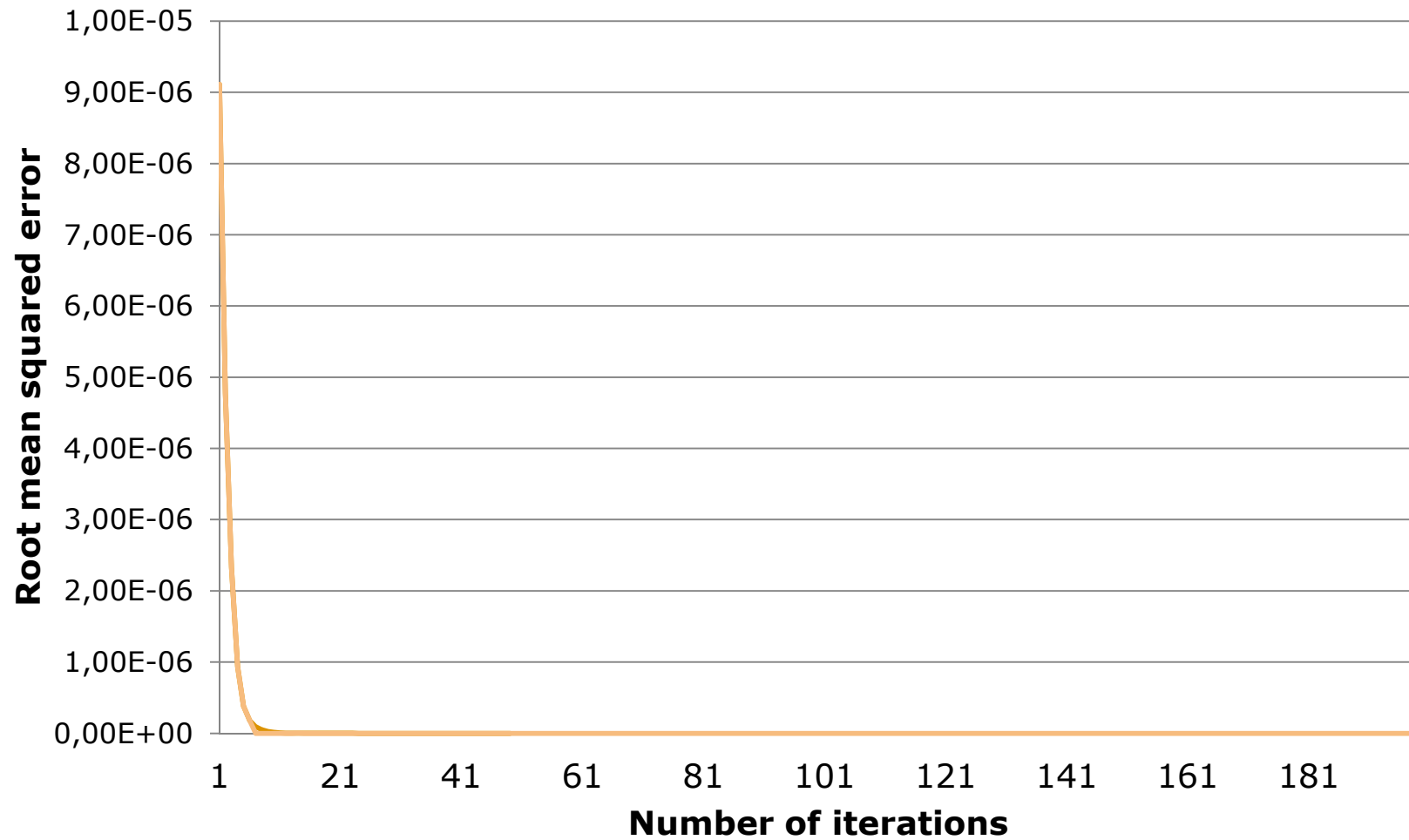
- Stratosphere
 - Explicit degree of parallelism
 - Hangs sometimes without error



1. PageRank on Stratosphere
2. Hadoop vs. Stratosphere
- 3. Evaluation**

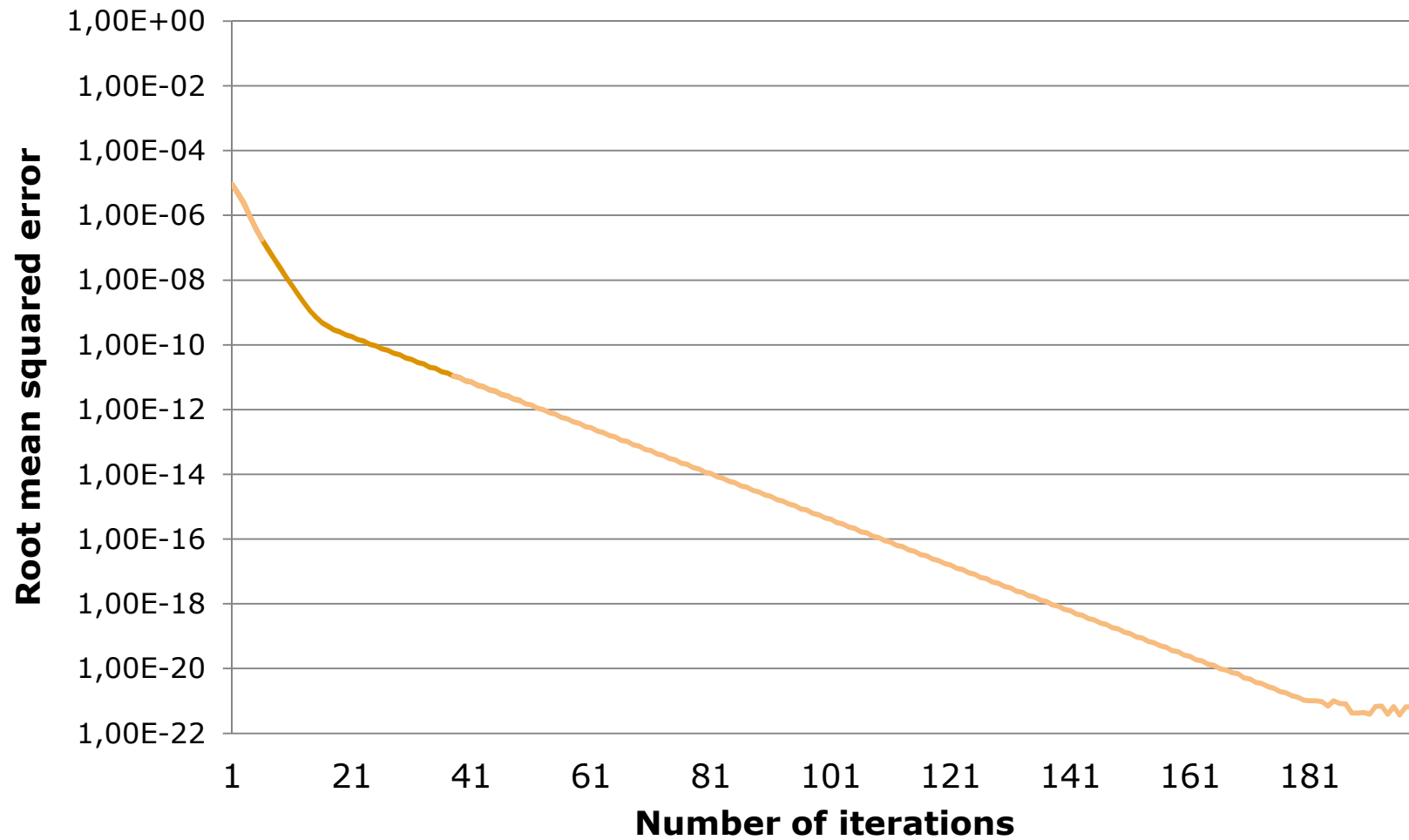
Convergence

35



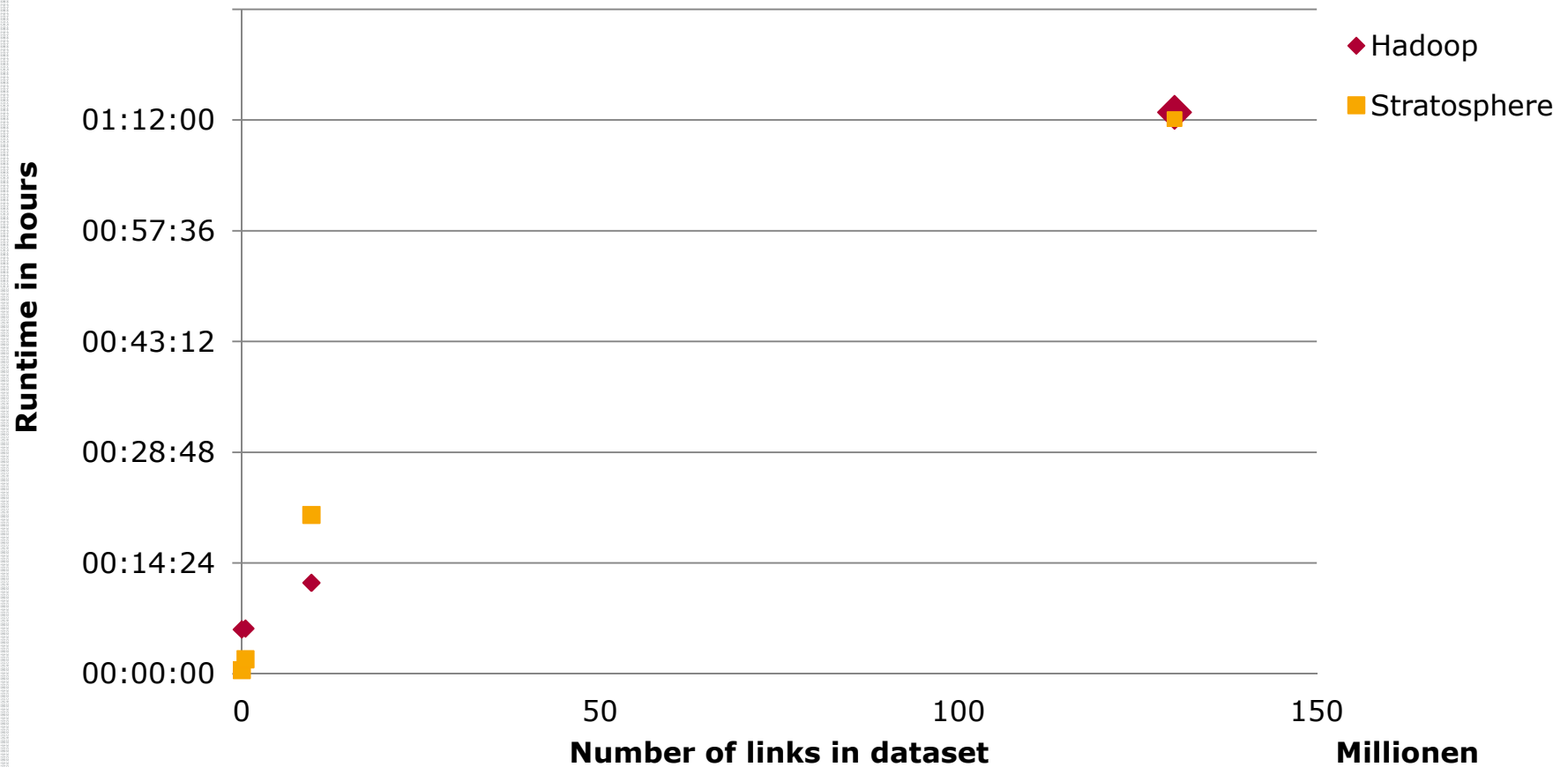
Convergence

36



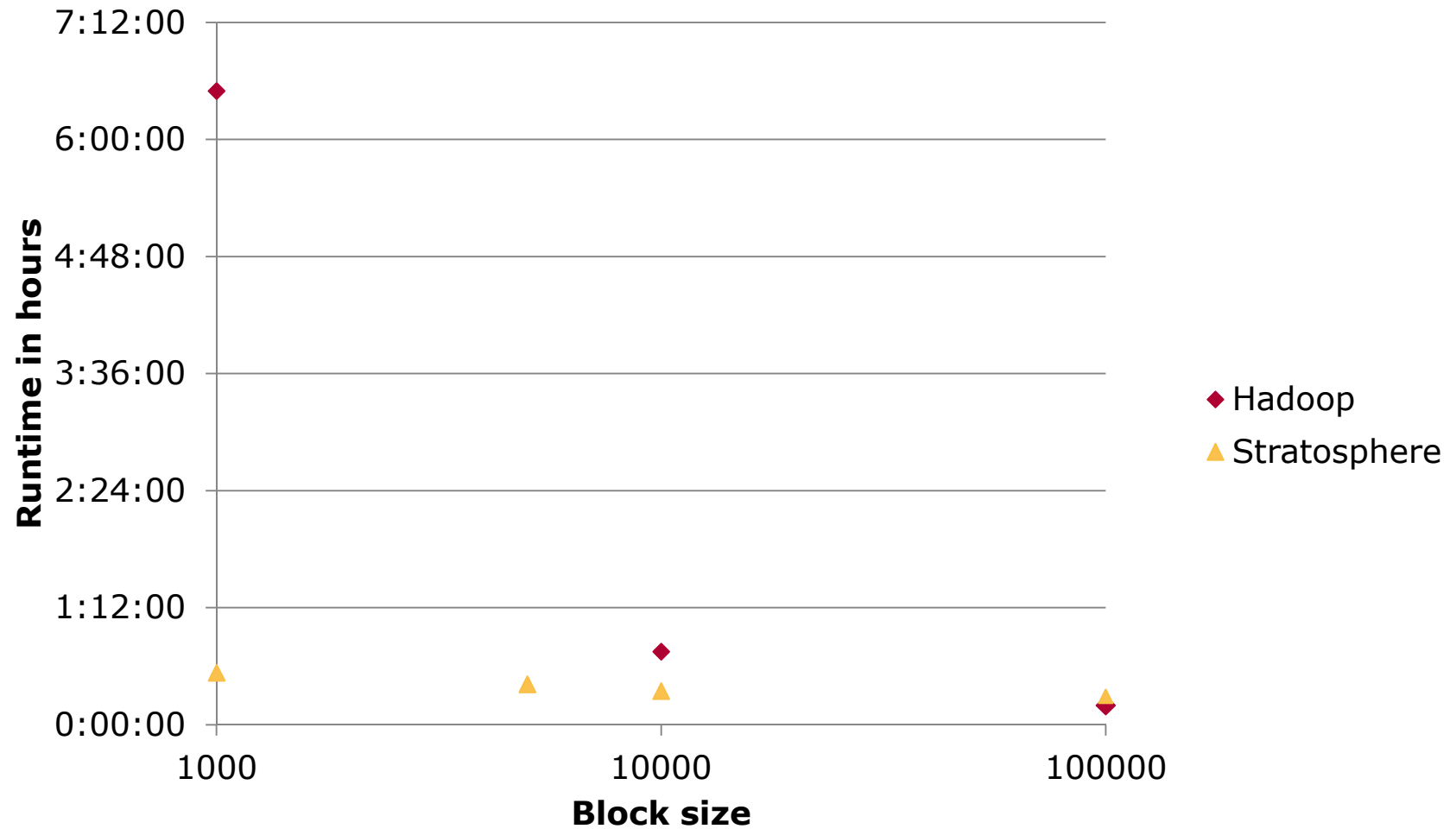
Runtime with different input sizes

37



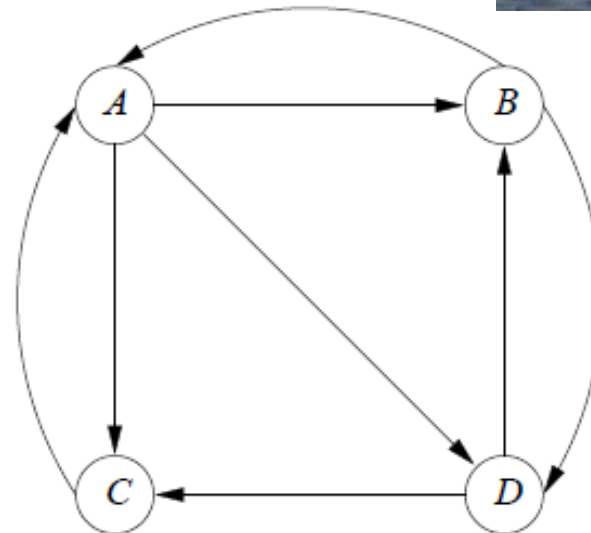
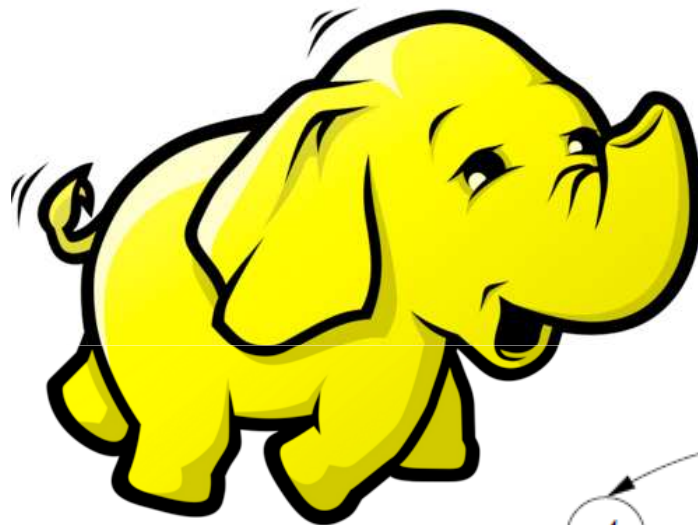
Runtime with different block sizes

38



Conclusion

39



Runtimes split up by job

