

Multi-pass Sorted Neighborhood Blocking with MapReduce

Lars Kolb, Andreas Thor, Erhard Rahm

Jens Hildebrandt, Jakob Zwiener

Agenda

2

1. Sorted Neighborhood Method
 - with Map Reduce
 - with Entity Replication
2. Multipass Sorted Neighborhood Method
3. Load Balancing
4. Benchmarks

Sorted Neighborhood Method

3

sorting key	artist_name	disc_title	Genre	tracks
	Sonny Terry	The Blues	Blues	18
	Fats Waller	Portrait	Jazz	17
	Blind Blake	Best Of	Blues	18
	Fats Domino	I'M Walking	Blues	18
	Chris Rea	Stony Road	Blues	17
	Jazz	Jazz	Jazz	20
	Acustica	Acustica	Blues	19
	Various	The Blues	Blues	17
	Kelis	Tasty	R+B	17

1. Calculate Sorting Key
 - Genre + tracks

Sorted Neighborhood Method

4

sorting key	artist_name	disc_title	Genre	tracks
Blues18	Sonny Terry	The Blues	Blues	18
Jazz17	Fats Waller	Portrait	Jazz	17
Blues18	Blind Blake	Best Of	Blues	18
Blues18	Fats Domino	I'M Walking	Blues	18
Blues17	Chris Rea	Stony Road	Blues	17
Jazz20	Jazz	Jazz	Jazz	20
Blues19	Acustica	Acustica	Blues	19
Blues17	Various	The Blues	Blues	17
R+B17	Kelis	Tasty	R+B	17

1. Calculate Sorting Key
 - Genre + tracks
2. Sort

Sorted Neighborhood Method

5

sorting key	artist_name	disc_title	Genre	tracks
Blues17	Chris Rea	Stony Road	Blues	17
Blues17	Various	The Blues	Blues	17
Blues18	Sonny Terry	The Blues	Blues	18
Blues18	Blind Blake	Best Of	Blues	18
Blues18	Fats Domino	I'M Walking	Blues	18
Blues19	Acustica	Acustica	Blues	19
Jazz17	Fats Waller	Portrait	Jazz	17
Jazz20	Jazz	Jazz	Jazz	20
R+B17	Kelis	Tasty	R+B	17

1. Calculate Sorting Key
 - Genre + tracks
2. Sort
3. Move a window over the data
 - Window size $w = 3$
 - Row count $n = 9$

Comparisons: $\mathcal{O}(n*w)$

Sorted Neighborhood with Map Reduce - Algorithm

6

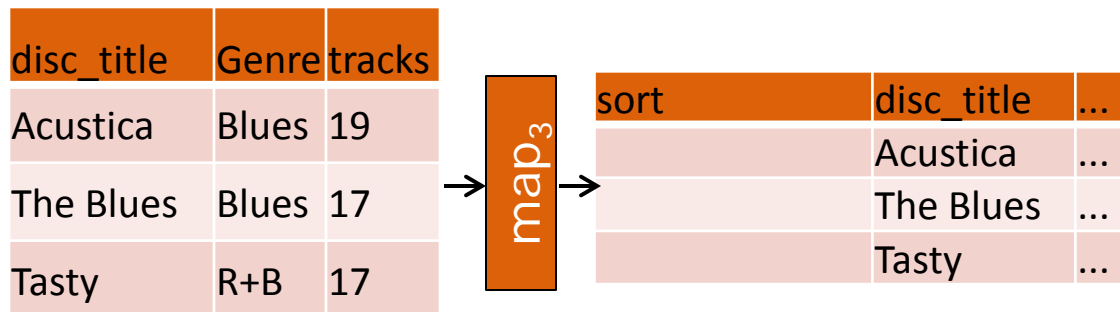
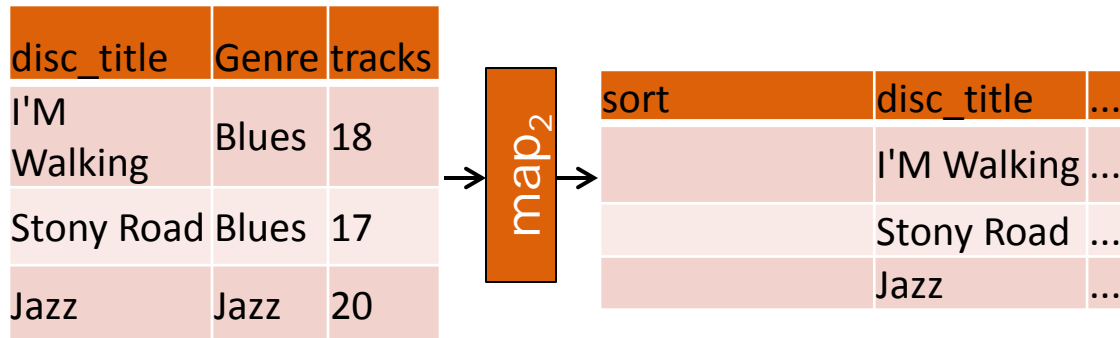
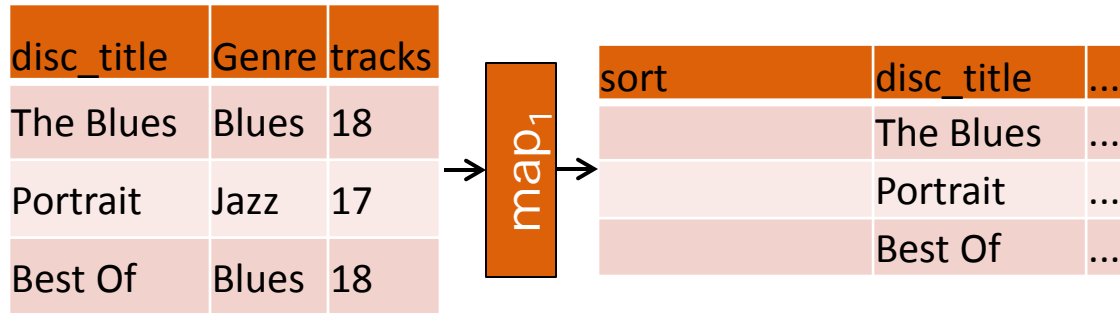
disc_title	Genre	tracks
The Blues	Blues	18
Portrait	Jazz	17
Best Of	Blues	18

disc_title	Genre	tracks
I'M Walking	Blues	18
Stony Road	Blues	17
Jazz	Jazz	20

disc_title	Genre	tracks
Acustica	Blues	19
The Blues	Blues	17
Tasty	R+B	17

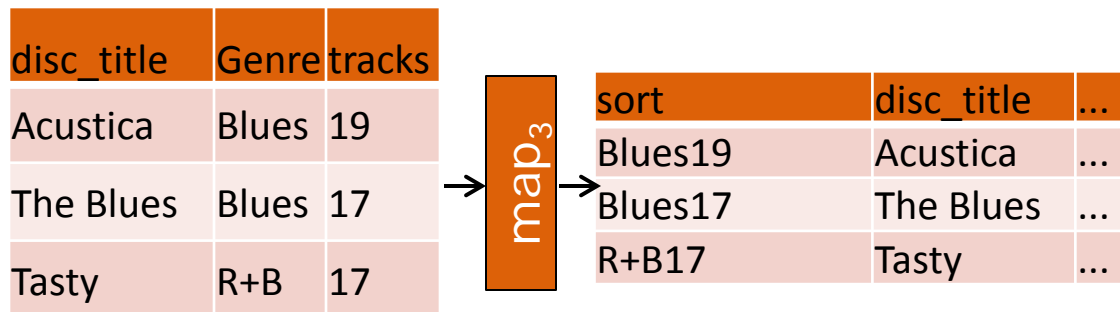
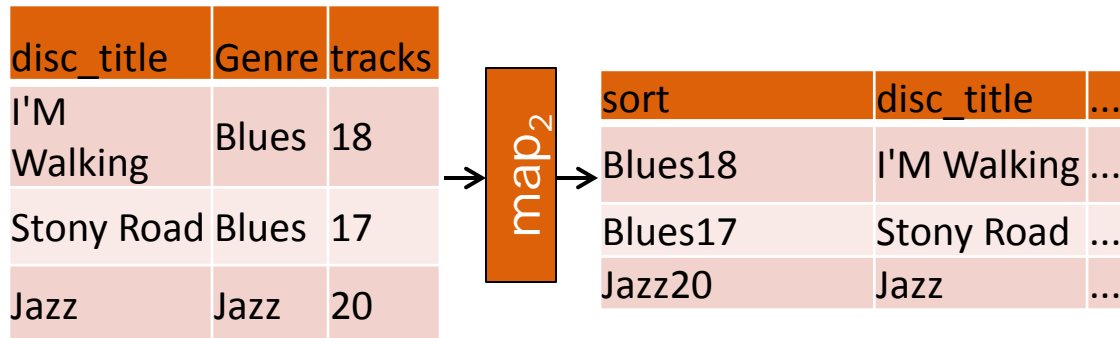
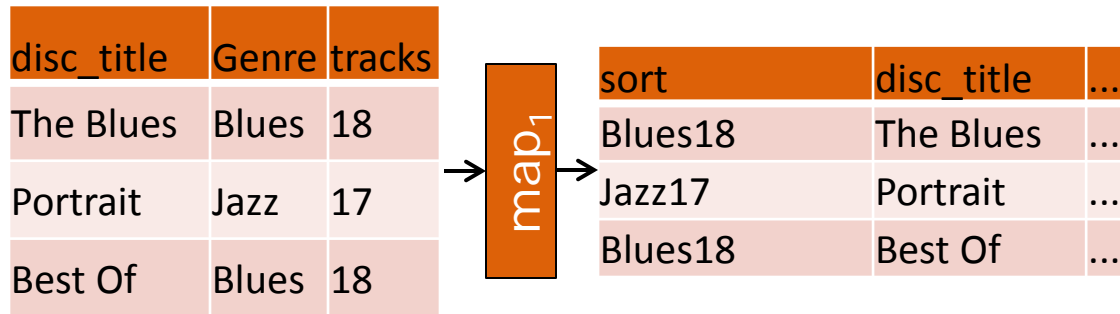
Sorted Neighborhood with Map Reduce - Algorithm

7



Sorted Neighborhood with Map Reduce - Algorithm

8

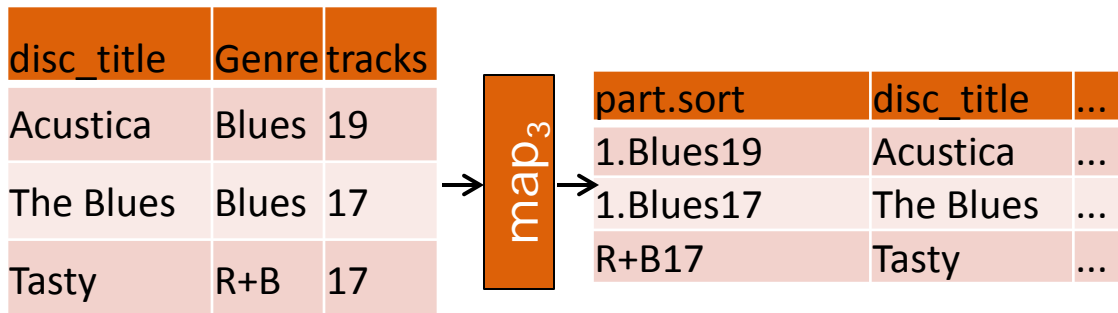
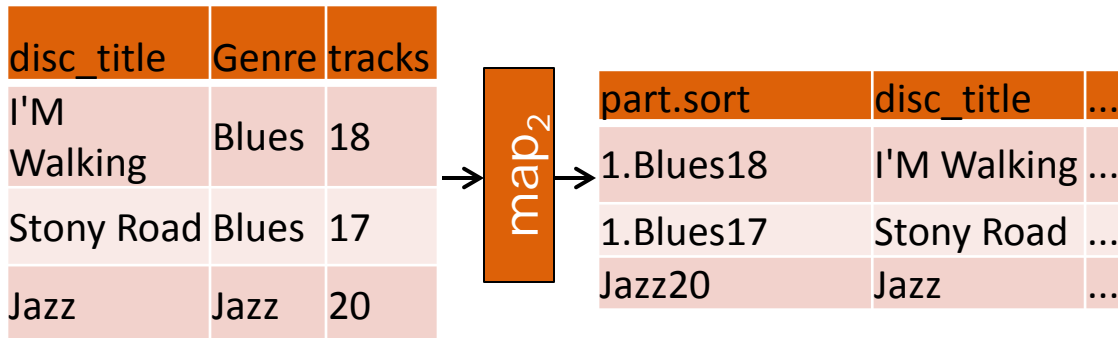
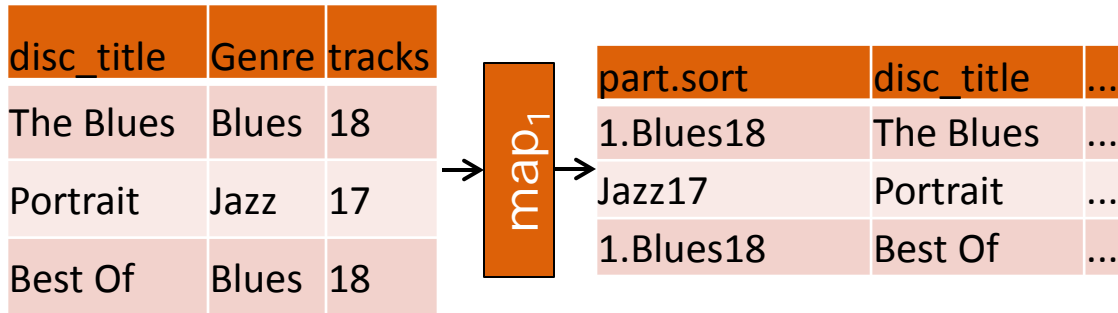


Map:

1. Calculate SortingKey: Genre+tracks

Sorted Neighborhood with Map Reduce - Algorithm

9



Map:

1. Calculate SortingKey:
Genre+tracks

2. Calculate Partition:

sorting key	partition
B...	1

Sorted Neighborhood with Map Reduce - Algorithm

10

disc_title	Genre	tracks
The Blues	Blues	18
Portrait	Jazz	17
Best Of	Blues	18

→ **map₁** →

part.sort	disc_title	...
1.Blues18	The Blues	...
2.Jazz17	Portrait	...
1.Blues18	Best Of	...

disc_title	Genre	tracks
I'M Walking	Blues	18
Stony Road	Blues	17
Jazz	Jazz	20

→ **map₂** →

part.sort	disc_title	...
1.Blues18	I'M Walking	...
1.Blues17	Stony Road	...
2.Jazz20	Jazz	...

disc_title	Genre	tracks
Acustica	Blues	19
The Blues	Blues	17
Tasty	R+B	17

→ **map₃** →

part.sort	disc_title	...
1.Blues19	Acustica	...
1.Blues17	The Blues	...
R+B17	Tasty	...

Map:

1. Calculate SortingKey: Genre+tracks

2. Calculate Partition:

sorting key	partition
B...	1
J...	2

Sorted Neighborhood with Map Reduce - Algorithm

11

disc_title	Genre	tracks
The Blues	Blues	18
Portrait	Jazz	17
Best Of	Blues	18

→ **map₁** →

part.sort	disc_title	...
1.Blues18	The Blues	...
2.Jazz17	Portrait	...
1.Blues18	Best Of	...

disc_title	Genre	tracks
I'M Walking	Blues	18
Stony Road	Blues	17
Jazz	Jazz	20

→ **map₂** →

part.sort	disc_title	...
1.Blues18	I'M Walking	...
1.Blues17	Stony Road	...
2.Jazz20	Jazz	...

disc_title	Genre	tracks
Acustica	Blues	19
The Blues	Blues	17
Tasty	R+B	17

→ **map₃** →

part.sort	disc_title	...
1.Blues19	Acustica	...
1.Blues17	The Blues	...
2.R+B17	Tasty	...

Map:

1. Calculate SortingKey: Genre+tracks

2. Calculate Partition:

sorting key	partition
B...	1
J...	2
R...	2

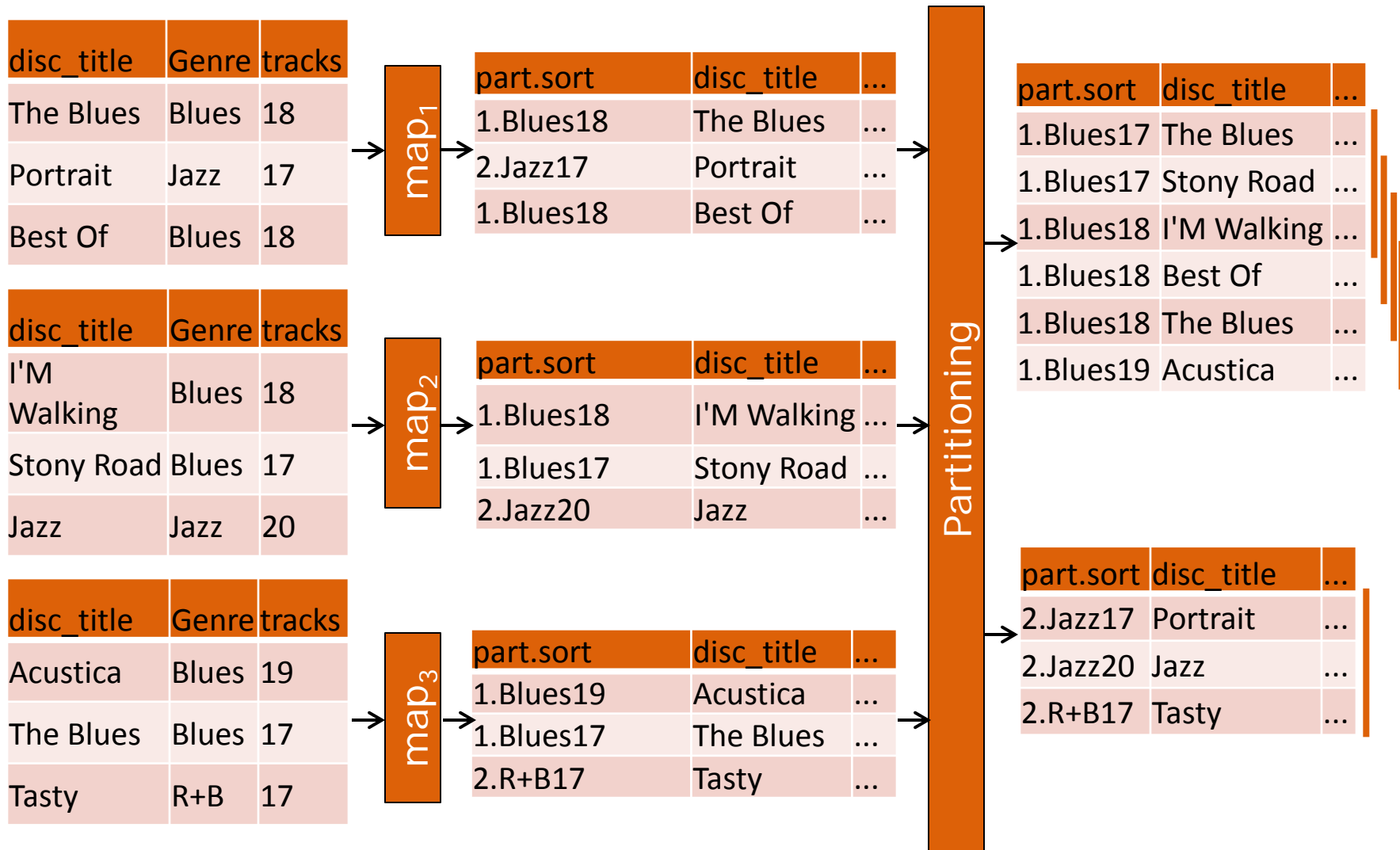
Sorted Neighborhood with Map Reduce - Algorithm

12



Sorted Neighborhood with Map Reduce - Algorithm

13



Sorted Neighborhood with Map Reduce - Limitations

14

- Neighboring sorting keys must be on the same reducer
 - own partition function
 - Self defined partitioning + sorting
 - Internal load balancing does not work anymore
- Boundary entities
 - Sliding window cannot compare entities that are assigned to different reduce nodes
 - Solution: data replication

reduce₁

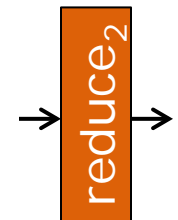
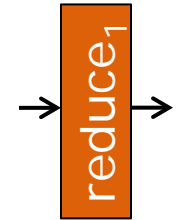
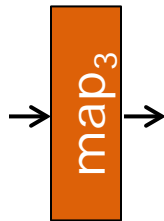
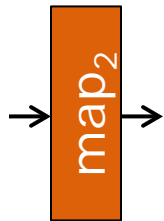
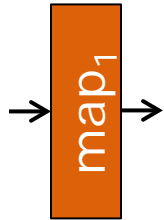
part.sort	disc_title	...
1.Blues17	The Blues	...
1.Blues17	Stony Road	...
1.Blues18	I'M Walking	...
1.Blues18	Best Of	...
1.Blues18	The Blues	...
1.Blues19	Acustica	...

reduce₂

part.sort	disc_title	...
2.Jazz17	Portrait	...
2.Jazz20	Jazz	...
2.R+B17	Tasty	...

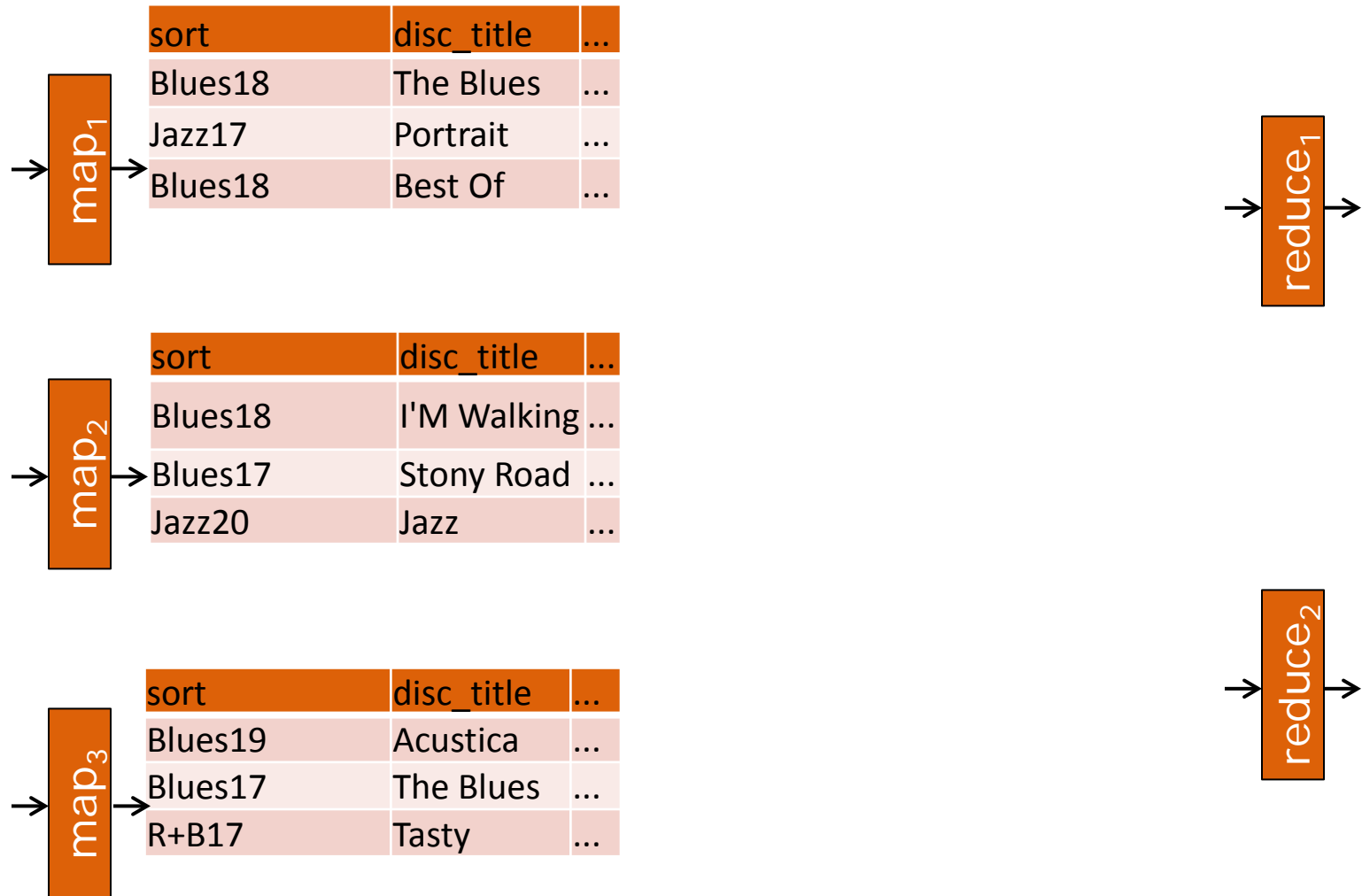
Sorted Neighborhood with Entity Replication

15



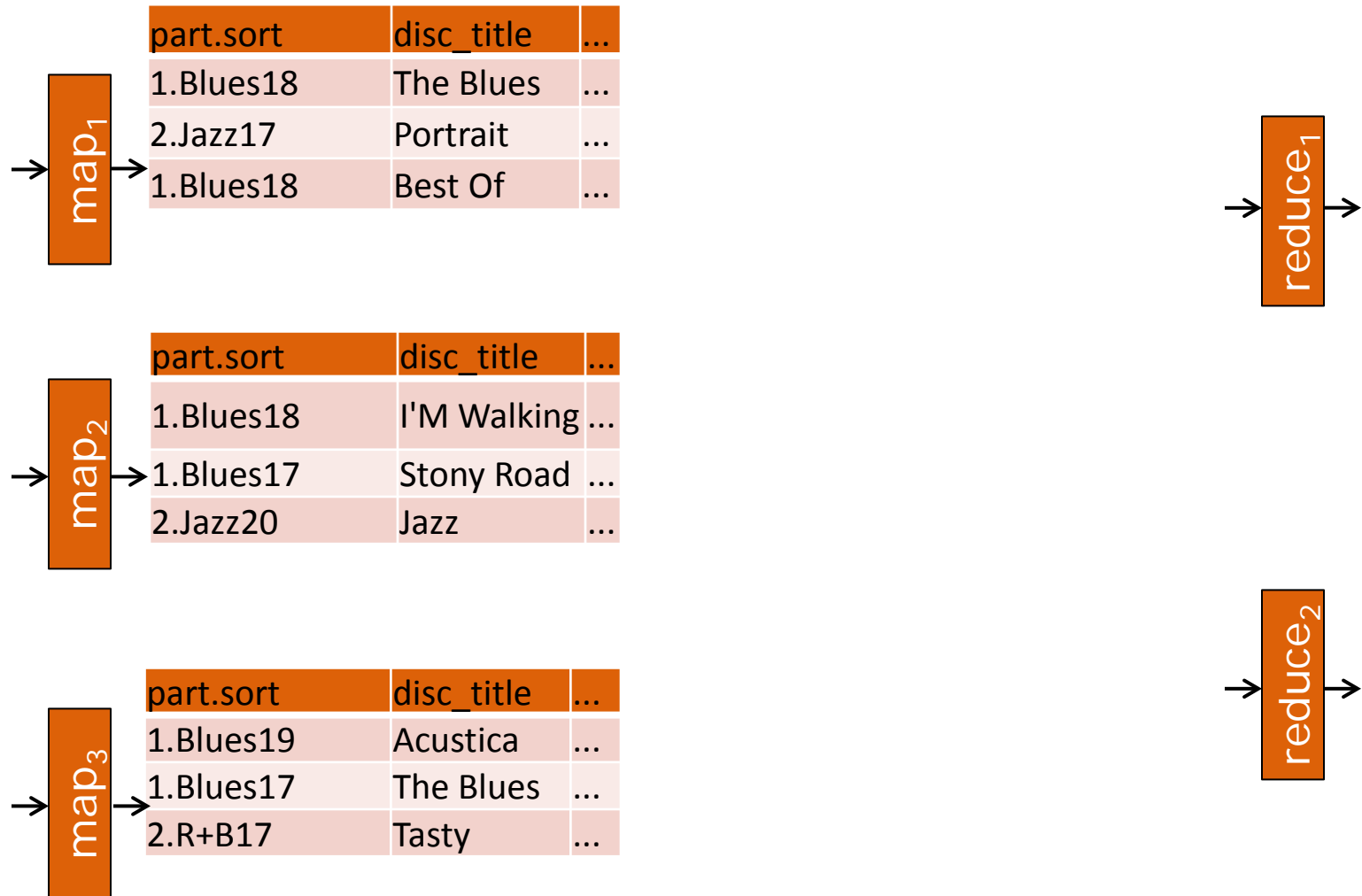
Sorted Neighborhood with Entity Replication

16



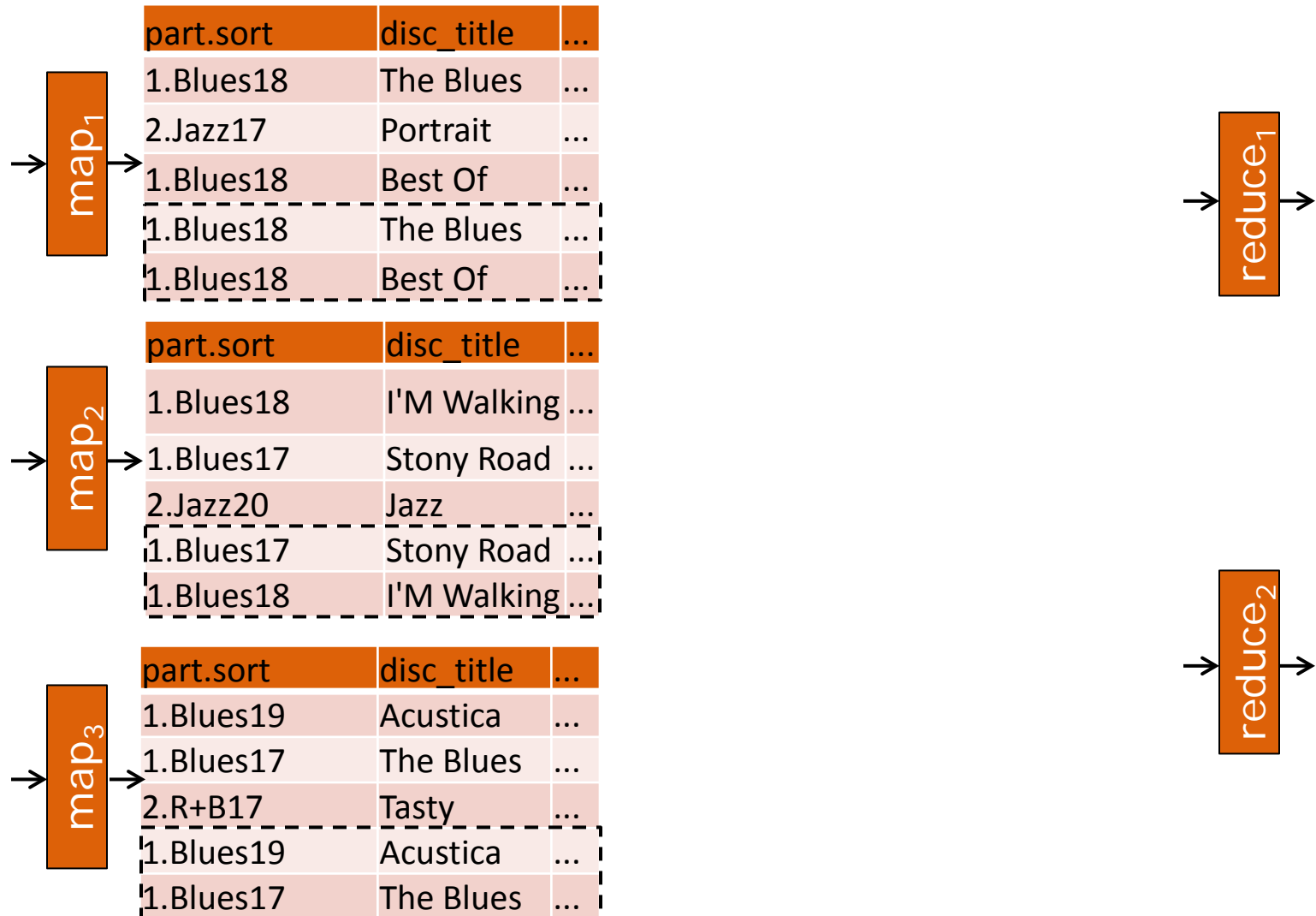
Sorted Neighborhood with Entity Replication

17



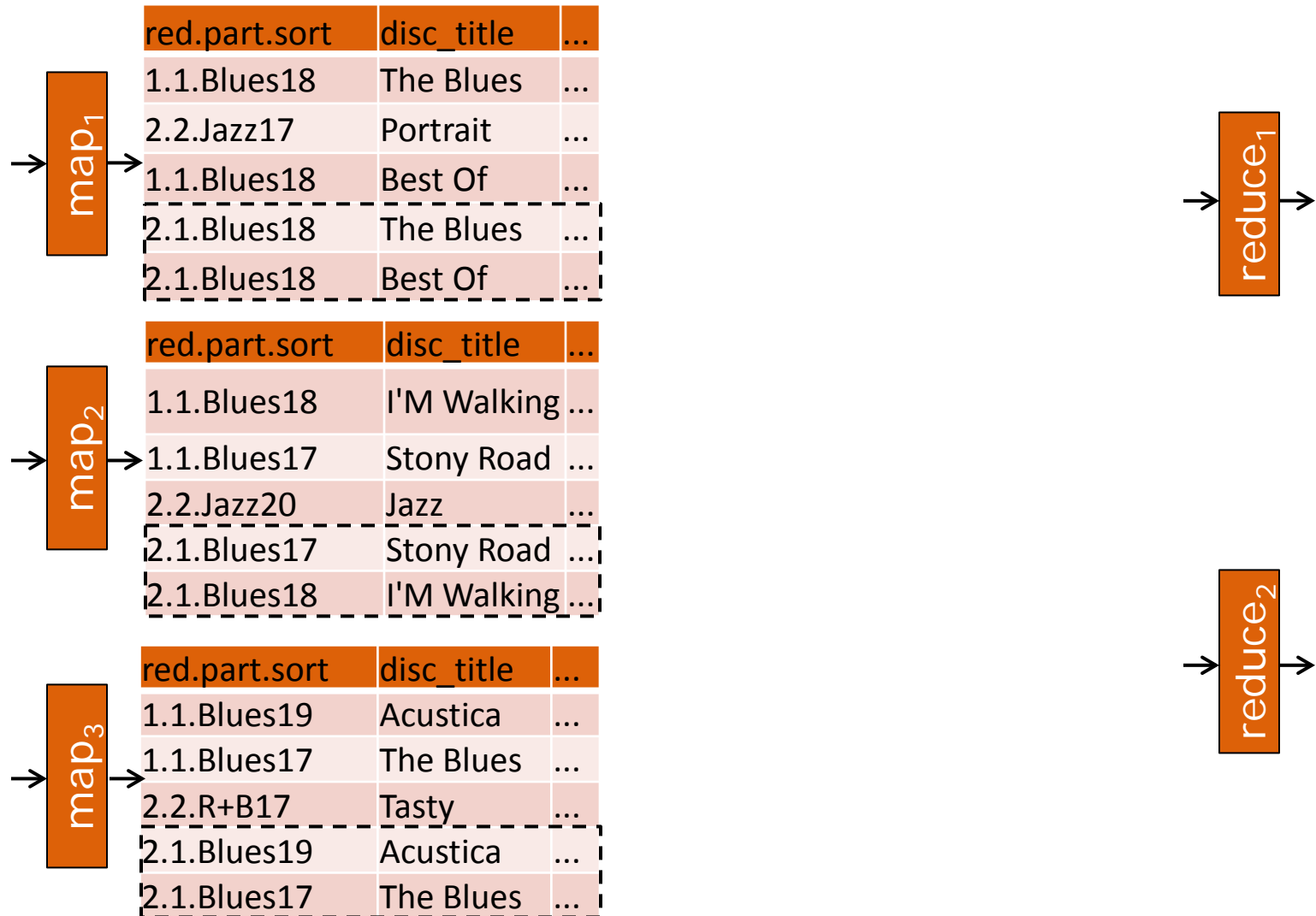
Sorted Neighborhood with Entity Replication

18



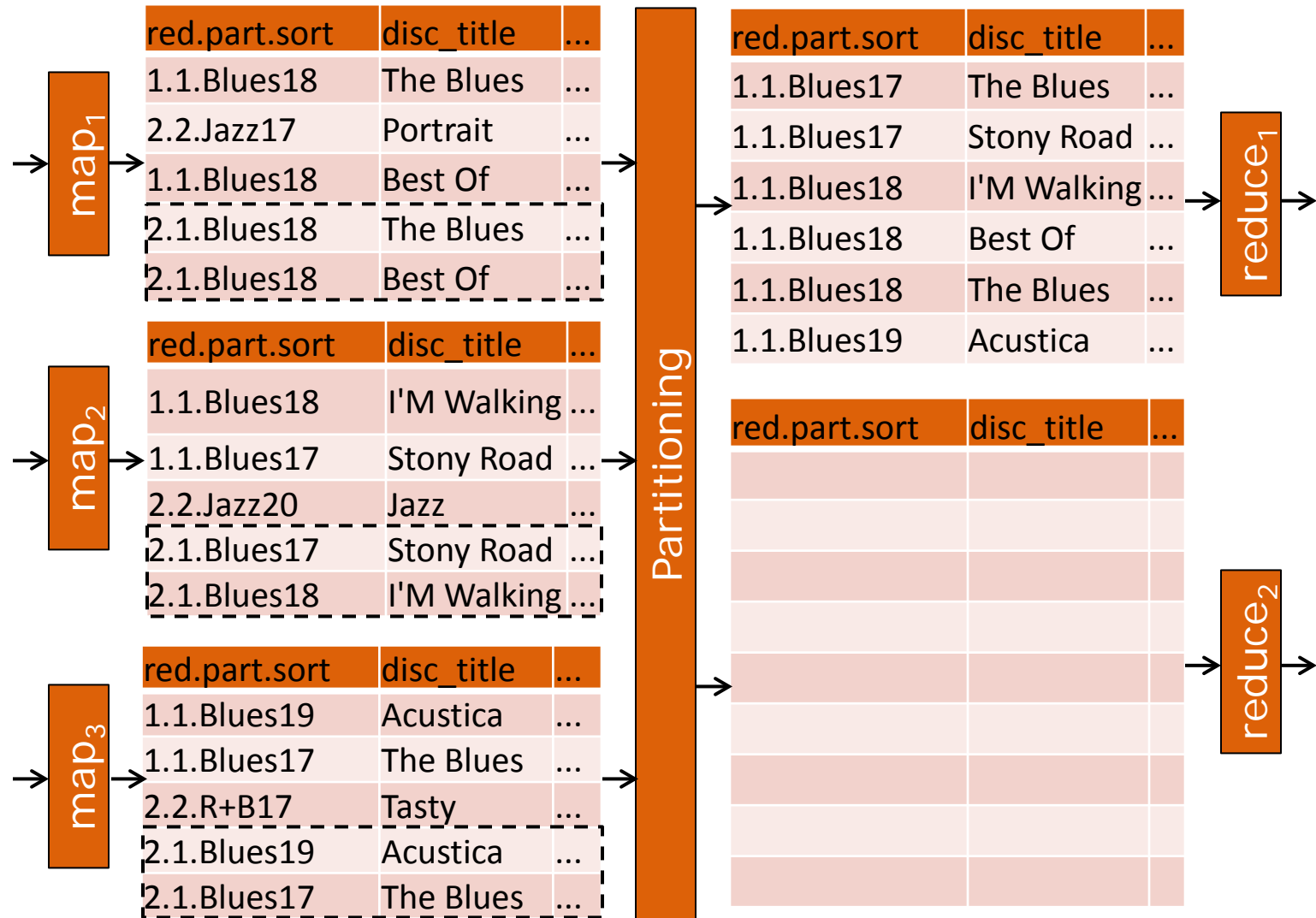
Sorted Neighborhood with Entity Replication

19



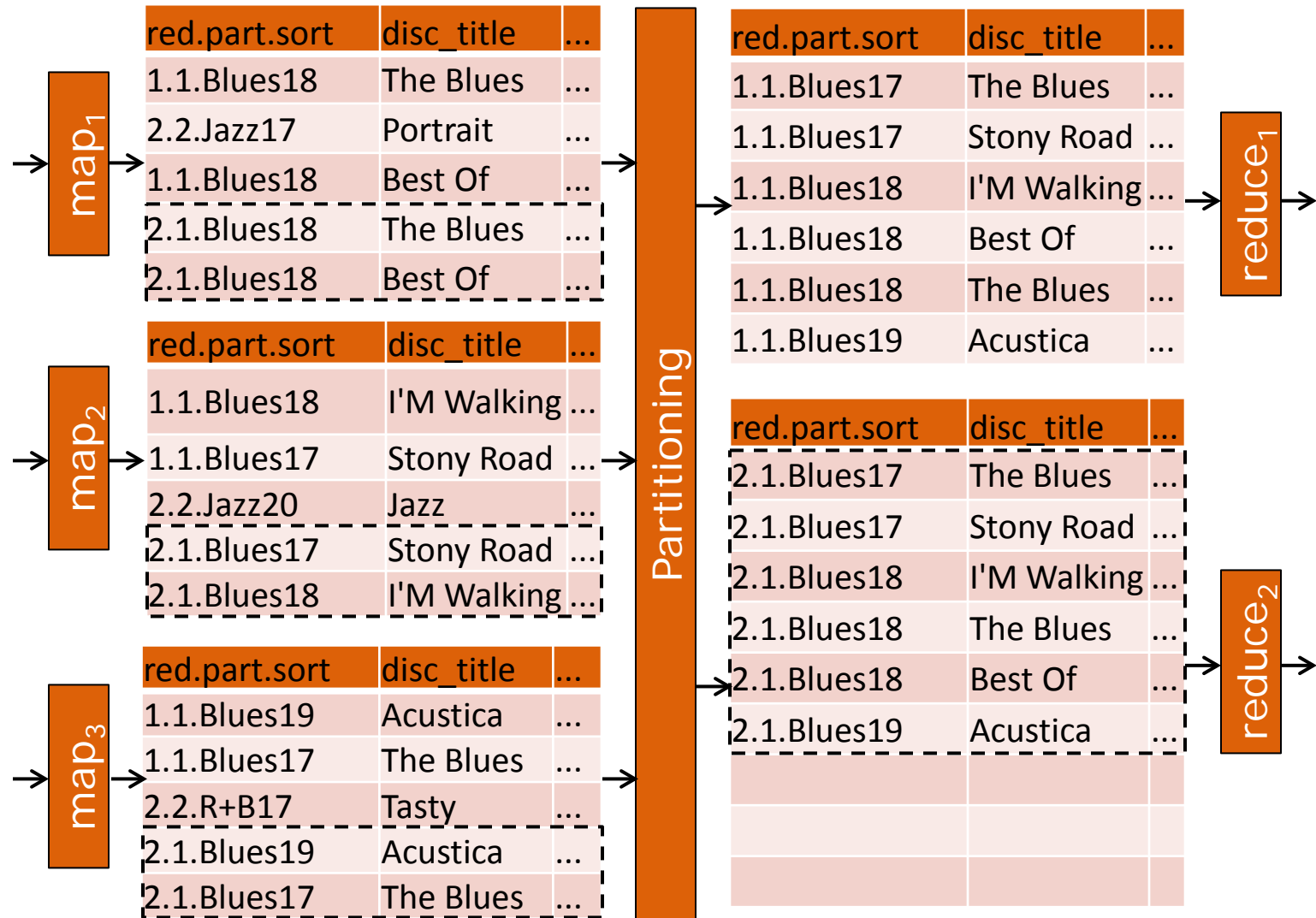
Sorted Neighborhood with Entity Replication

20



Sorted Neighborhood with Entity Replication

21



Sorted Neighborhood with Entity Replication

22



Challenges in Sorted Neighborhood on Map Reduce

24

- Sorted Neighborhood with Map Reduce ✓
- Multipass in one Map Reduce
- Load Balancing for Nodes

Multipass Sorted Neighborhood Method

25

disc_title	Genre	tracks
The Blues	Blues	18
Portrait	Jazz	17
Best Of	Blues	18

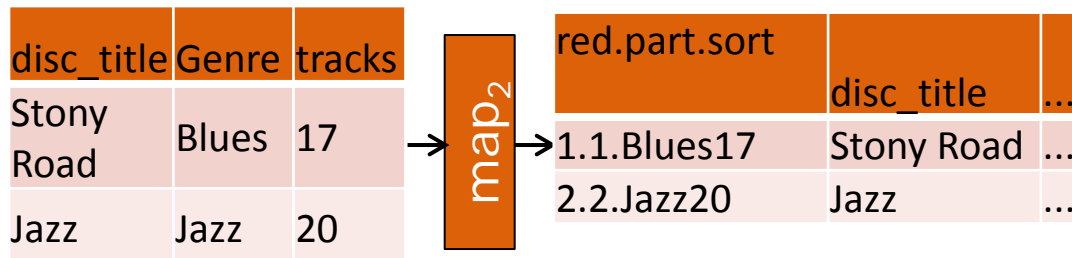
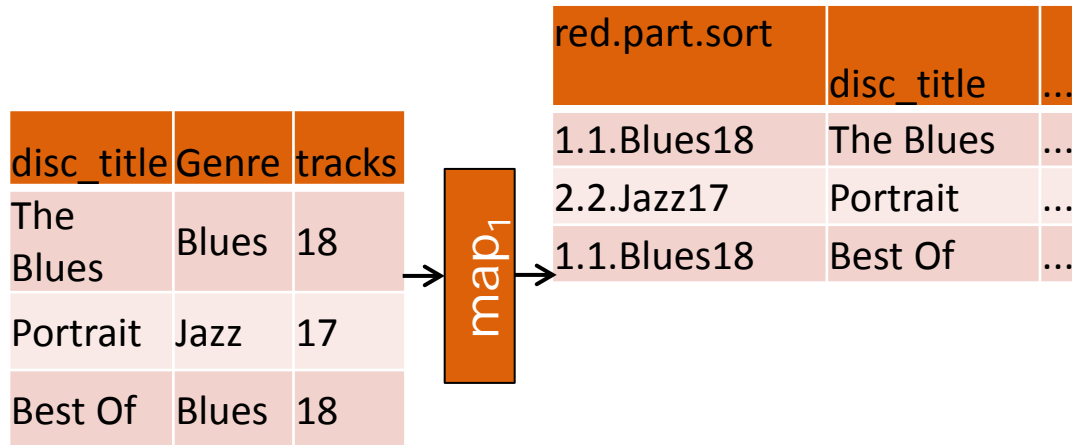
→ map₁

disc_title	Genre	tracks
Stony Road	Blues	17
Jazz	Jazz	20

→ map₂

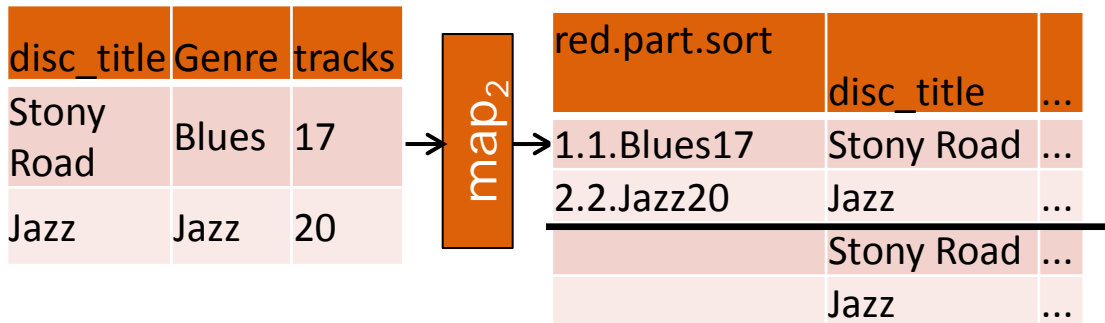
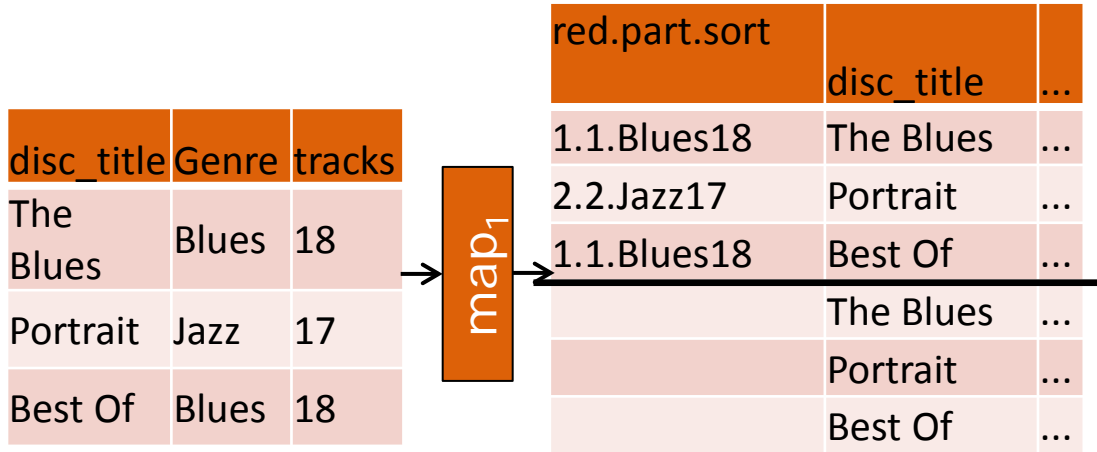
Multipass Sorted Neighborhood Method

26



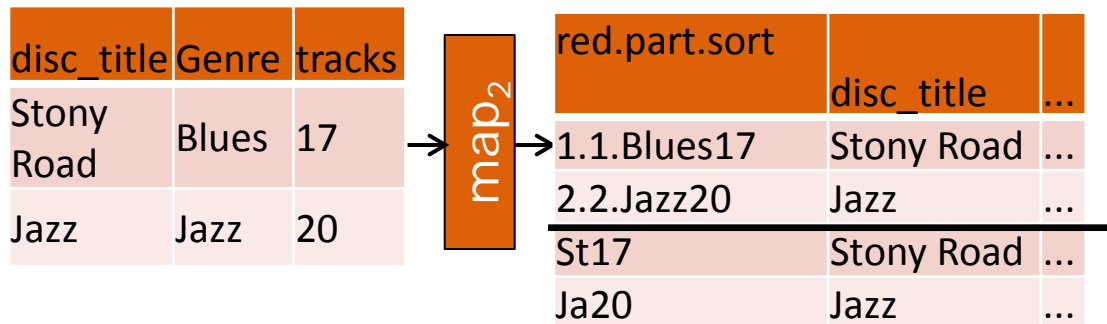
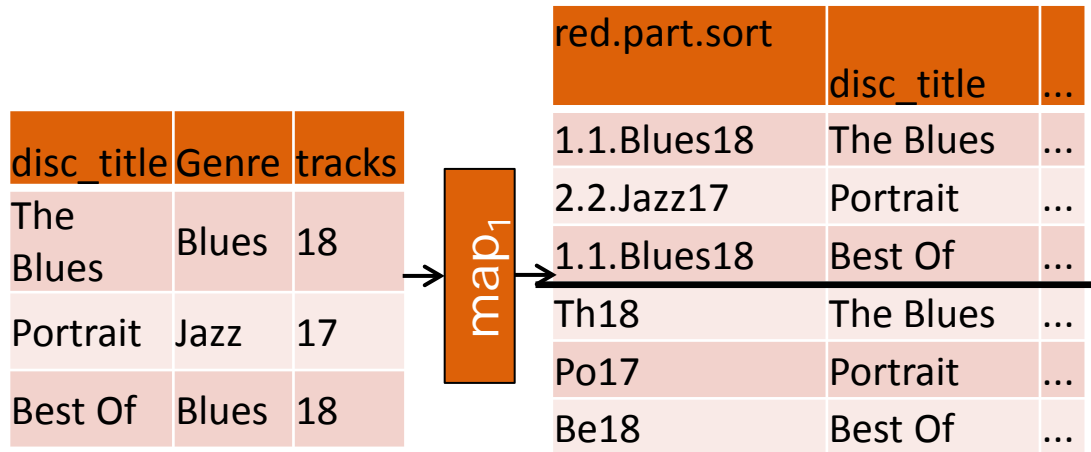
Multipass Sorted Neighborhood Method

27



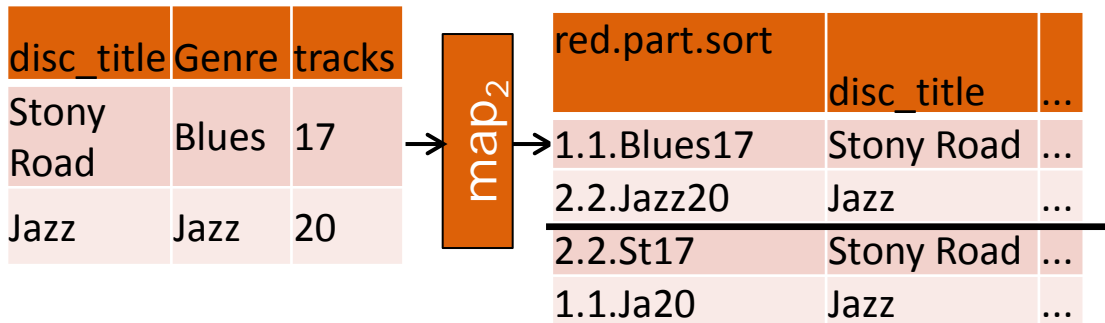
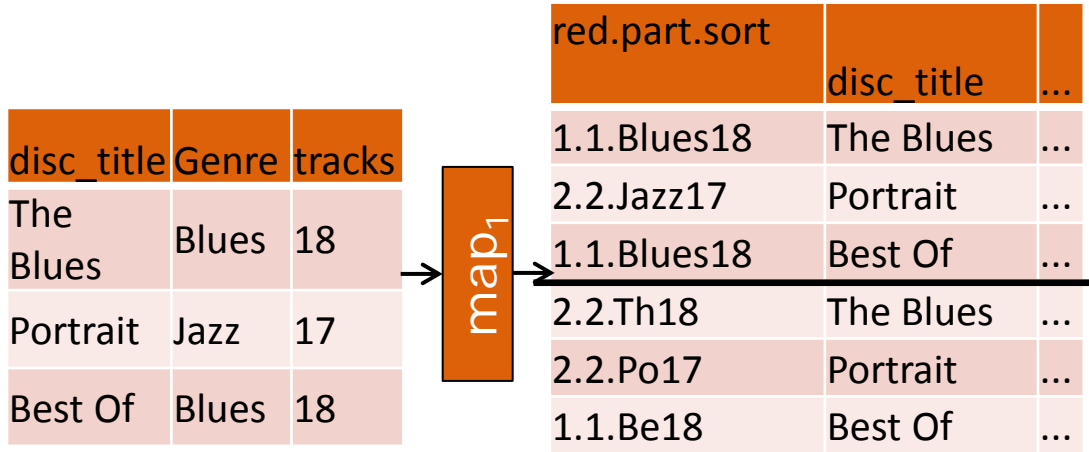
Multipass Sorted Neighborhood Method

28



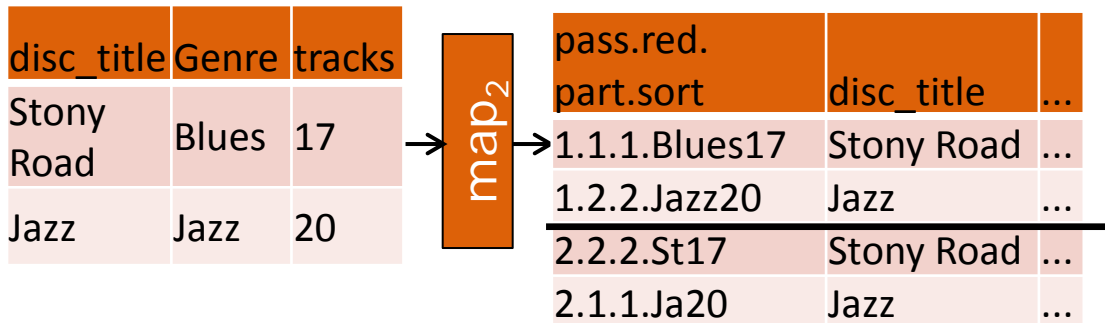
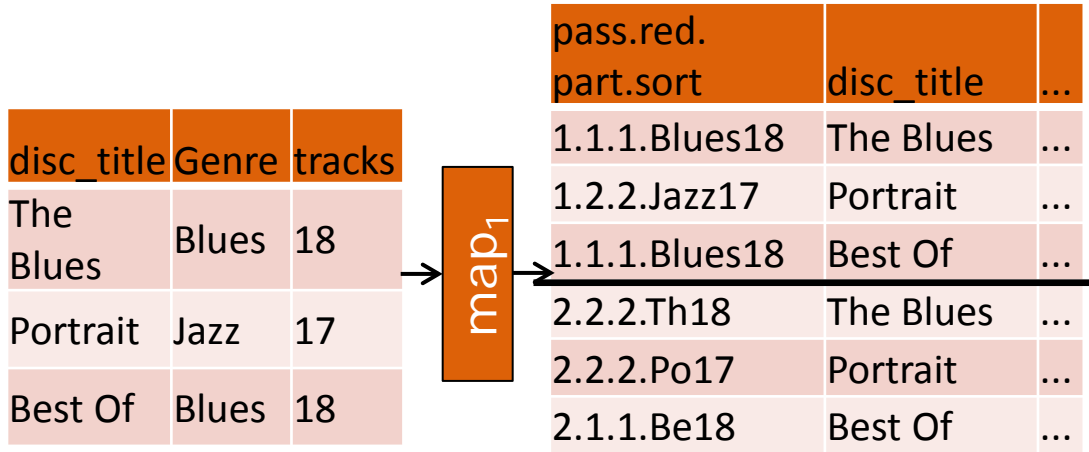
Multipass Sorted Neighborhood Method

29



Multipass Sorted Neighborhood Method

30



Challenges in Sorted Neighborhood on Map Reduce

33

- Sorted Neighborhood with Map Reduce ✓
- Multipass in one Map Reduce ✓
- Load Balancing for Nodes

Load Balancing

34

sort	disc_title	...
Blues18	The Blues	...
Jazz17	Portrait	...
Blues18	Best Of	...

sort	disc_title	...
Blues18	I'M Walking	...
Blues17	Stony Road	...
Jazz20	Jazz	...

sort	disc_title	...
Blues19	Acustica	...
Blues17	The Blues	...
R+B17	Tasty	...

sort	disc_title	...
Blues17	Stony Road	...
Blues17	The Blues	...
Blues18	The Blues	...
Blues18	Best Of	...
Blues18	Best Of	...
Blues18	I'M Walking	...

sort	disc_title	...
Blues18	I'M Walking	...
Blues19	Acustica	...
Jazz17	Portrait	...
Jazz20	Jazz	...
R+B17	Tasty	...

Load Balancing

35

sort.mapN	disc_title	...
Blues18.1	The Blues	...
Jazz17.1	Portrait	...
Blues18.1	Best Of	...

sort.mapN	disc_title	...
Blues17.2	Stony Road	...
Blues17.3	The Blues	...
Blues18.1	The Blues	...
Blues18.1	Best Of	...

sort.mapN	disc_title	...
Blues18.2	I'M Walking	...
Blues17.2	Stony Road	...
Jazz20.2	Jazz	...

sort.mapN	disc_title	...
Blues18.2	I'M Walking	...
Blues19.3	Acustica	...
Jazz17.1	Portrait	...
Jazz20.2	Jazz	...
R+B17.3	Tasty	...

sort.mapN	disc_title	...
Blues19.3	Acustica	...
Blues17.3	The Blues	...
R+B17.3	Tasty	...

Load Balancing

36

sort.mapN.counter	disc_title	...
Blues18.1.1	The Blues	...
Jazz17.1.1	Portrait	...
Blues18.1.2	Best Of	...

sort.mapN.counter	disc_title	...
Blues17.2.1	Stony Road	...
Blues17.3.1	The Blues	...
Blues18.1.1	The Blues	...
Blues18.1.2	Best Of	...

sort.mapN.counter	disc_title	...
Blues18.2.1	I'M Walking	...
Blues17.2.1	Stony Road	...
Jazz20.2.1	Jazz	...

sort.mapN.counter	disc_title	...
Blues18.2.1	I'M Walking	...
Blues19.3.1	Acustica	...
Jazz17.1.1	Portrait	...
Jazz20.2.1	Jazz	...
R+B17.3.1	Tasty	...

sort.mapN.counter	disc_title	...
Blues19.3.1	Acustica	...
Blues17.3.1	The Blues	...
R+B17.3.1	Tasty	...

Load Balancing

37

sortKey	MapN:	1	2	3
Blues17		0	1	1
Blues18		2	1	0
Blues19		0	0	1
Jazz17		1	0	0
Jazz20		0	1	0
R+B17		0	0	1

Blues18.2.1

part.sort	disc_title	...

part.sort	disc_title	...

Load Balancing

38

sortKey	MapN:	1	2	3
Blues17		0	1	1
Blues18		2	1	0
Blues19		0	0	1
Jazz17		1	0	0
Jazz20		0	1	0
R+B17		0	0	1

Blues18.2.1

part.sort	disc_title	...

part.sort	disc_title	...
2.Blues18	I'M Walking	...

Load Balancing

39

sortKey	MapN:	1	2	3
Blues17		0	1	1
Blues18		2	1	0
Blues19		0	0	1
Jazz17		1	0	0
Jazz20		0	1	0
R+B17		0	0	1

Blues18.1.1

part.sort	disc_title	...

part.sort	disc_title	...
2.Blues18	I'M Walking	...

Load Balancing

40

sortKey	MapN:	1	2	3
Blues17		0	1	1
Blues18		2	1	0
Blues19		0	0	1
Jazz17		1	0	0
Jazz20		0	1	0
R+B17		0	0	1

Blues18.1.1

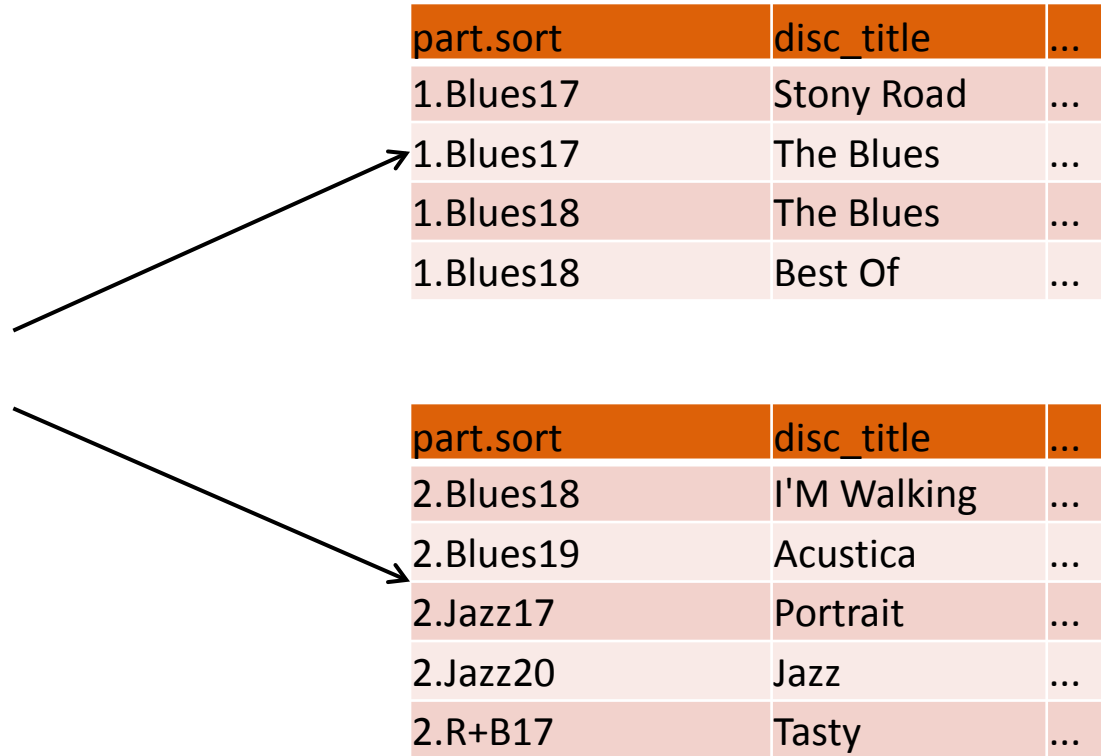
part.sort	disc_title	...
1.Blues18	The Blues	...

part.sort	disc_title	...
2.Blues18	I'M Walking	...

Load Balancing

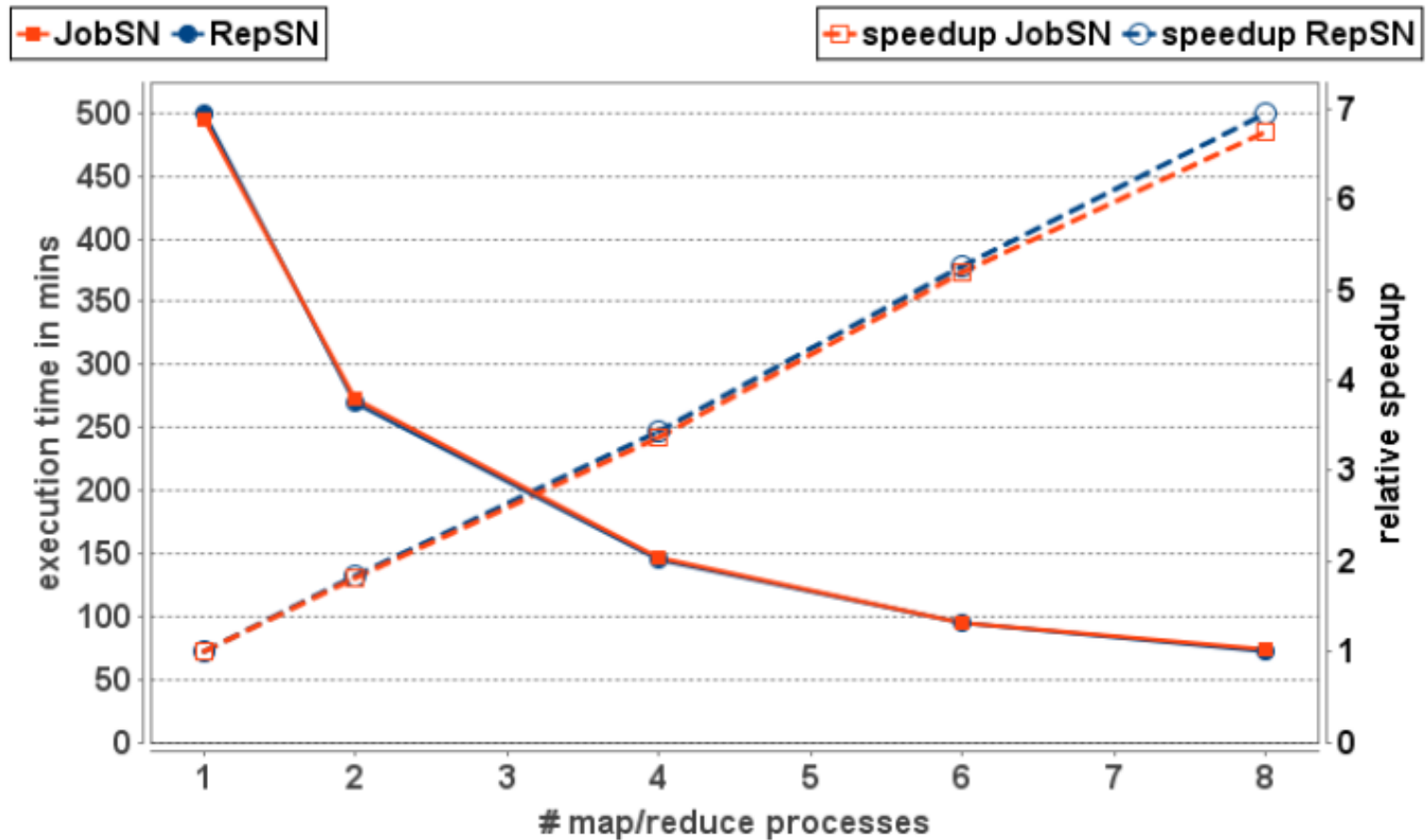
41

sortKey	MapN:	1	2	3
Blues17		0	1	1
Blues18		2	1	0
Blues19		0	0	1
Jazz17		1	0	0
Jazz20		0	1	0
R+B17		0	0	1



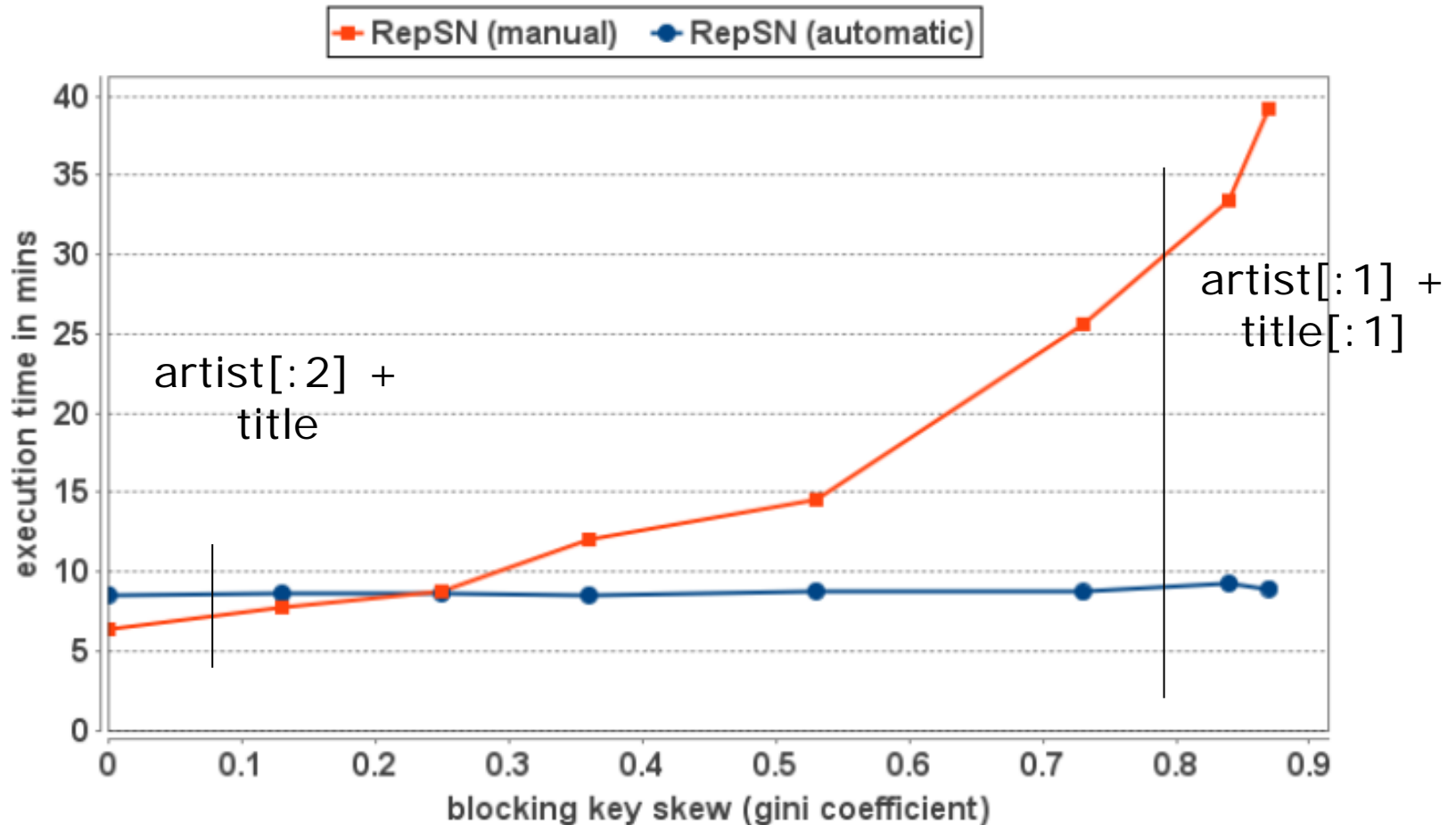
Benchmarks

43



Benchmarks

45



Summary

46

1. Sorted Neighborhood Method
 - with Map Reduce
 - with Entity Replication
2. Multipass Sorted Neighborhood Method
3. Load Balancing
4. Benchmarks