

Collaborative Filtering

Scalable Data Analysis Algorithms

Claudia Lehmann, Andrina Mascher

Outline

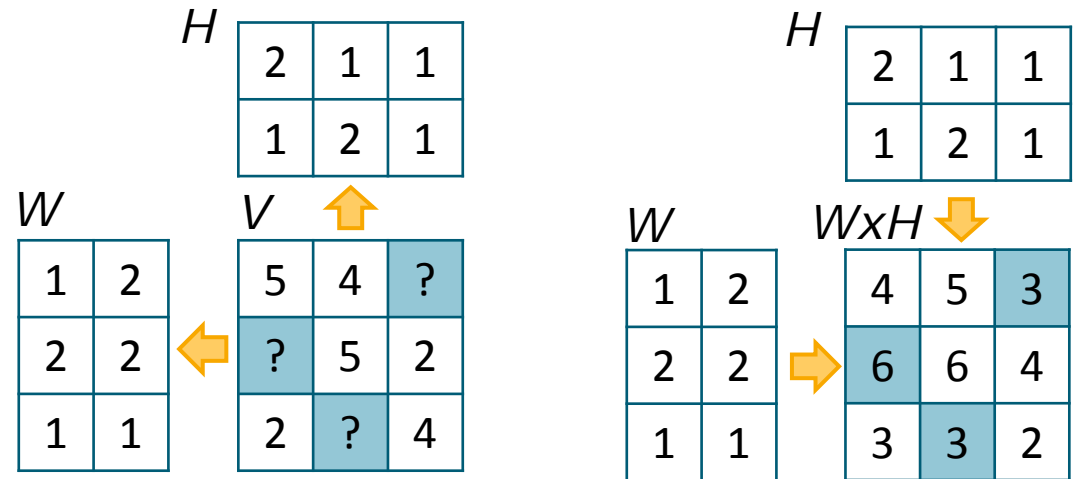
2

1. Retrospection
2. Stratosphere Plans
3. Comparison with Hadoop
4. Evaluation
5. Outlook

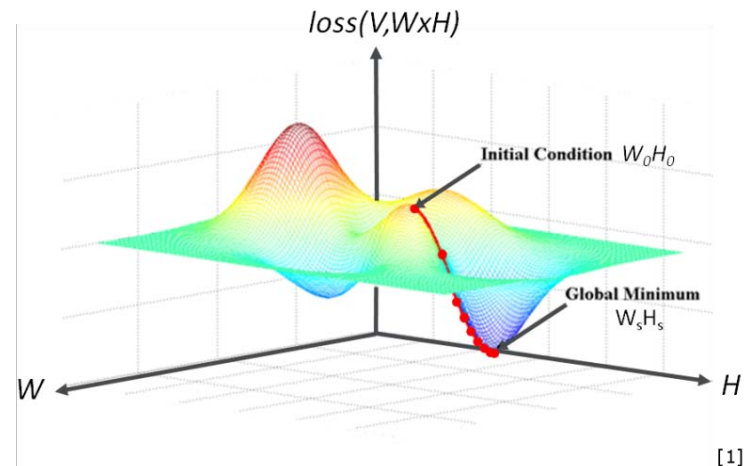
Retrospection

3

**Matrix factorization:
Find optimal factors**



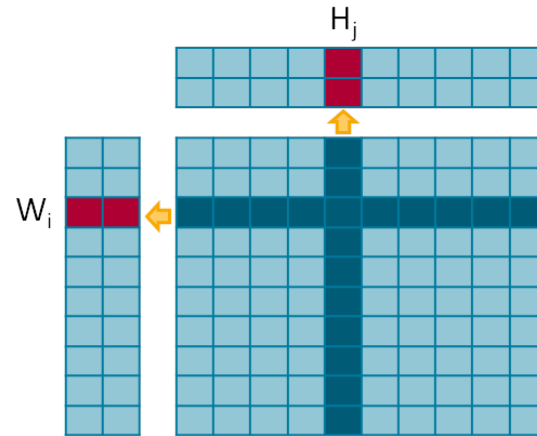
Stochastic Gradient Descent



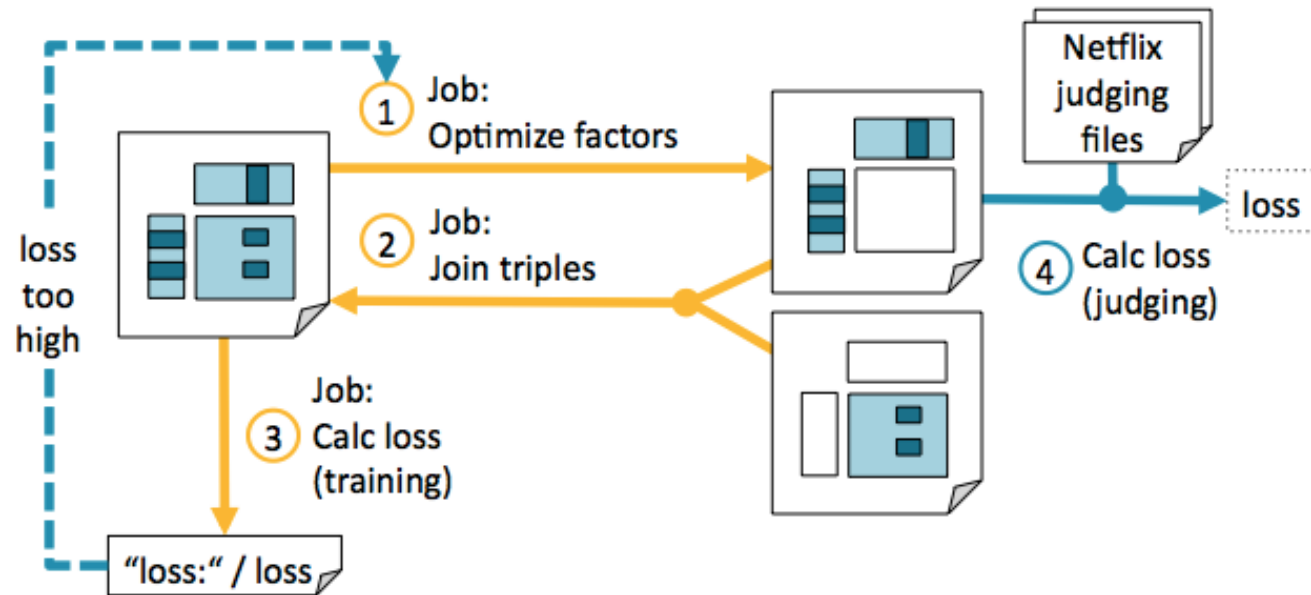
Retrospection

4

Result of SGD step



Hadoop jobs with judging loss



Stratosphere Plans

Multiple Iterations

5

```
for each starting point:  
  while loss too high:  
    MapReduce optimize factors  
    MapReduce join triples  
    MapReduce calculate loss (training)  
  calculate loss (judging)
```

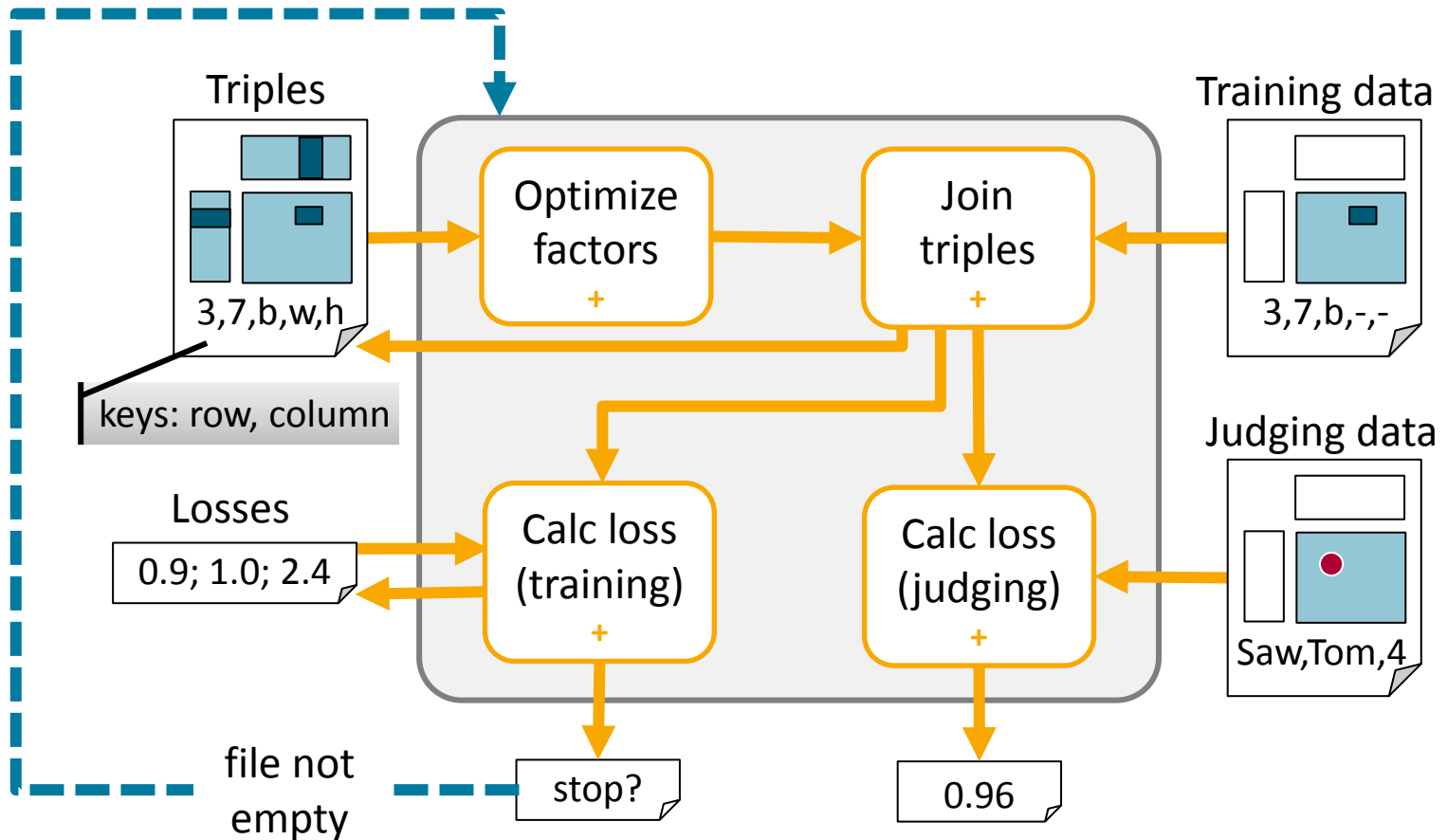


```
for each starting point:  
  while stop file not empty:  
    MapMapReduce optimize factors  
    CrossMatch join triples  
    MapReduceCross calculate loss (training) and decide  
    MapMatchReduce calculate loss (judging)
```

Stratosphere Plans

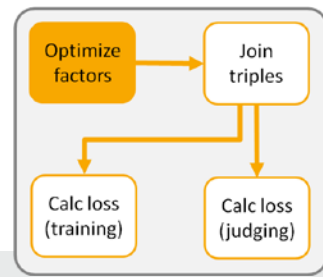
Subplans in One Iteration

6

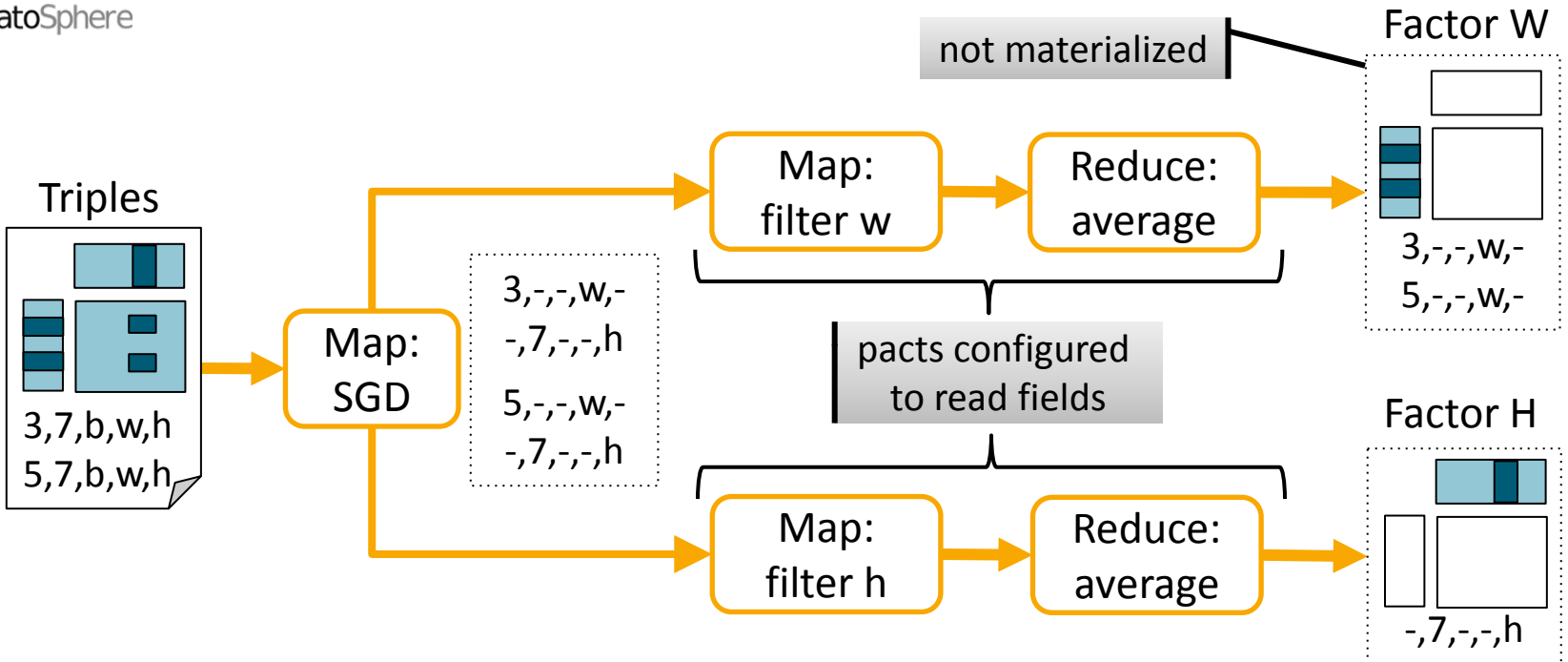
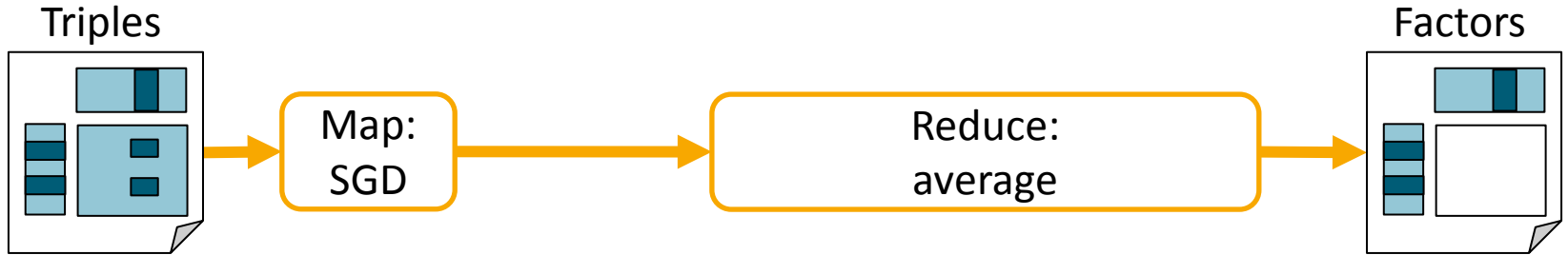


Stratosphere Plans

Optimize Factors

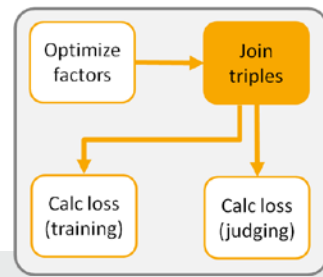


7

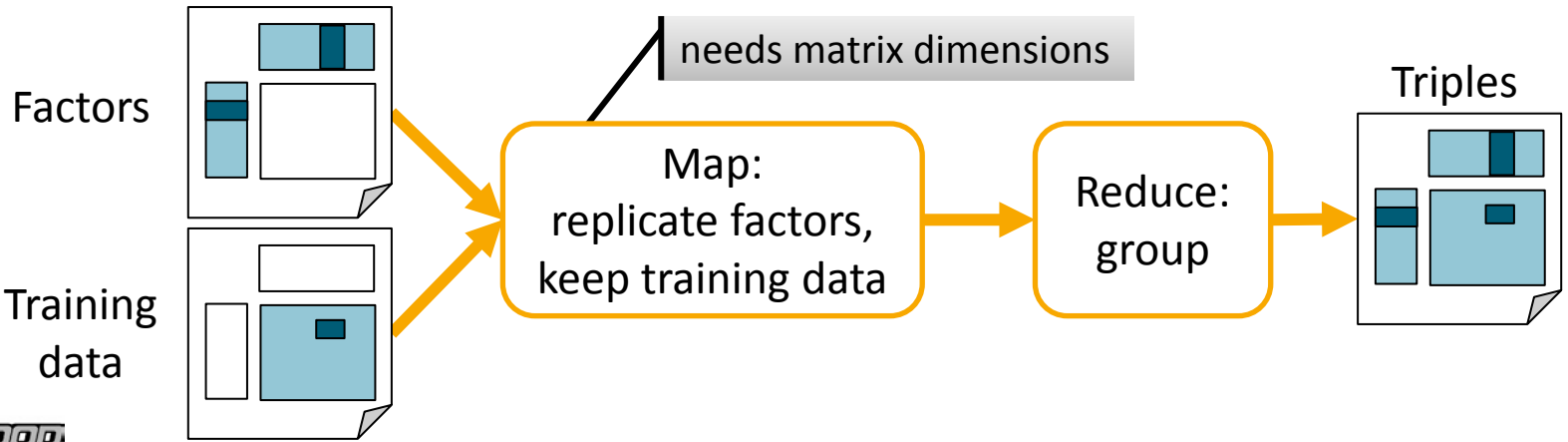


Stratosphere Plans

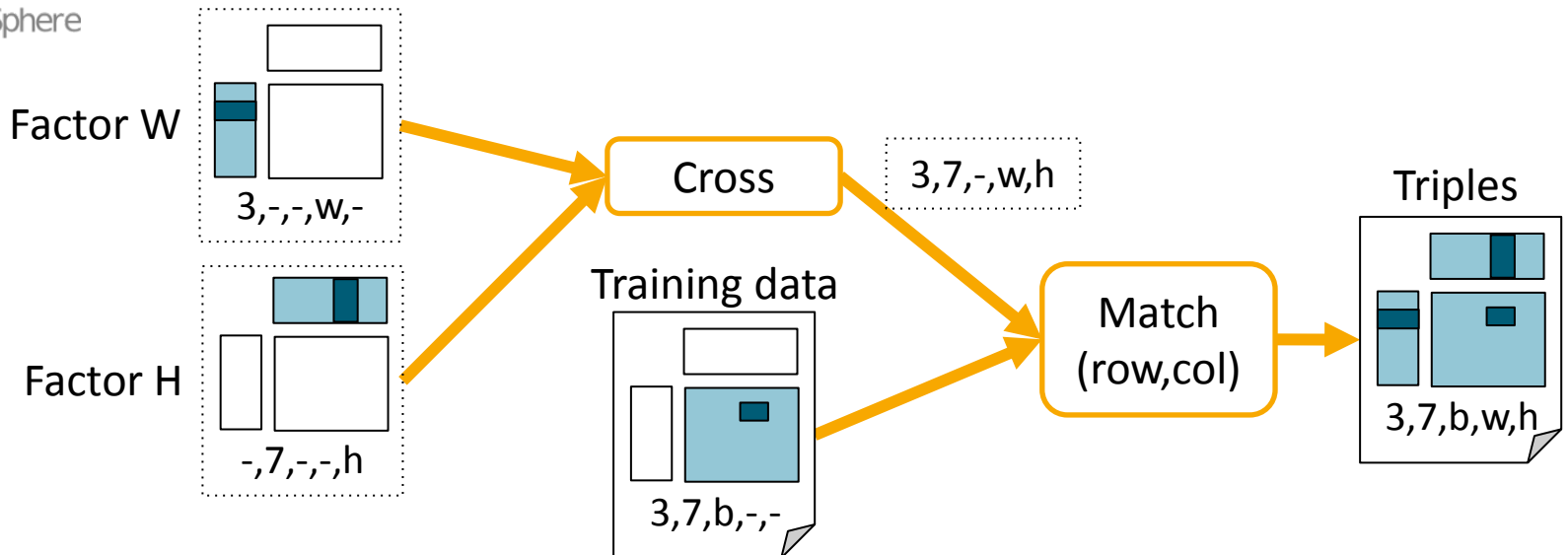
Join Triples (1/2)



8

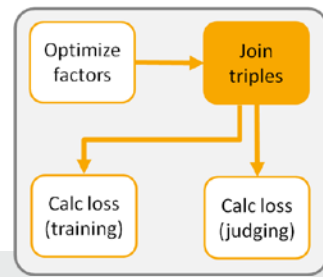


StratoSphere

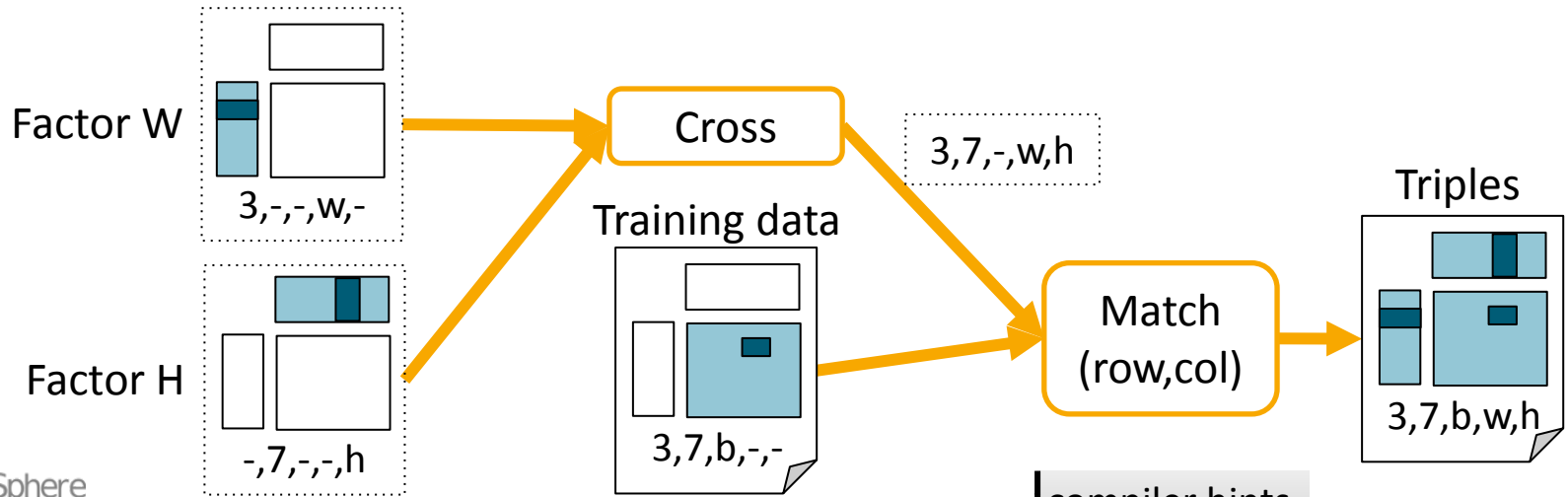


Stratosphere Plans

Join Triples (2/2)



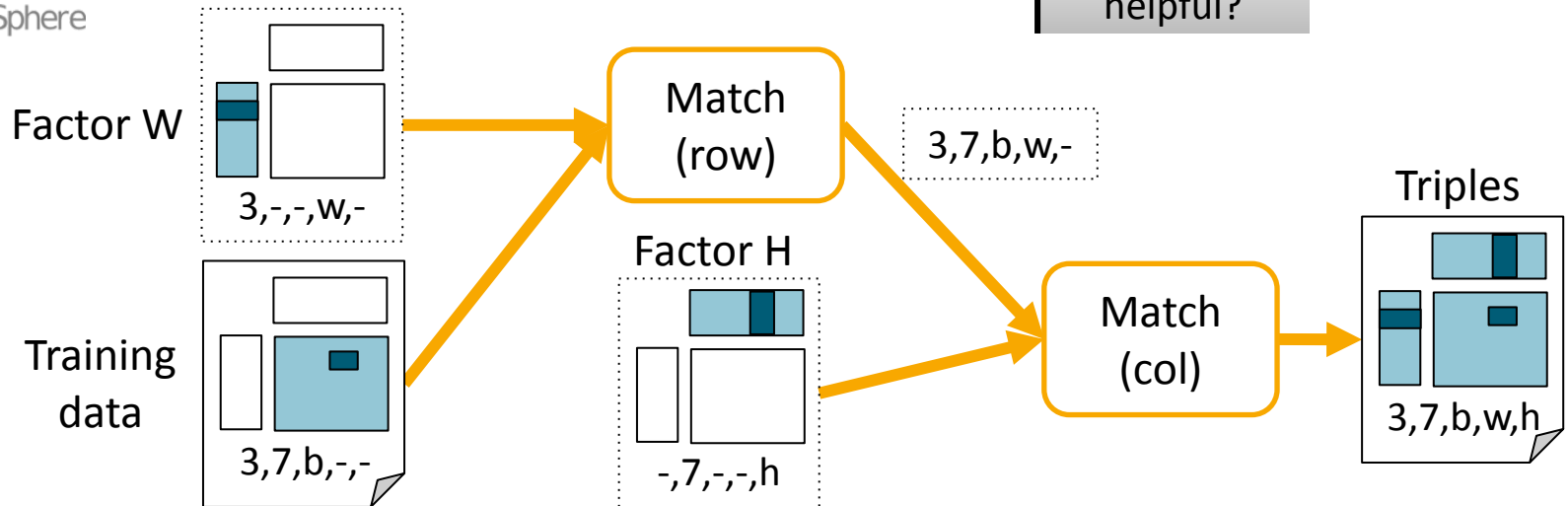
9



compiler hints helpful?

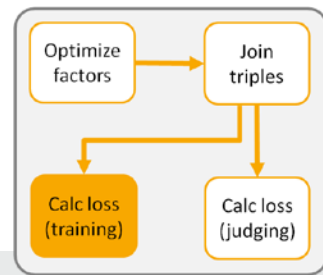
StratoSphere

StratoSphere

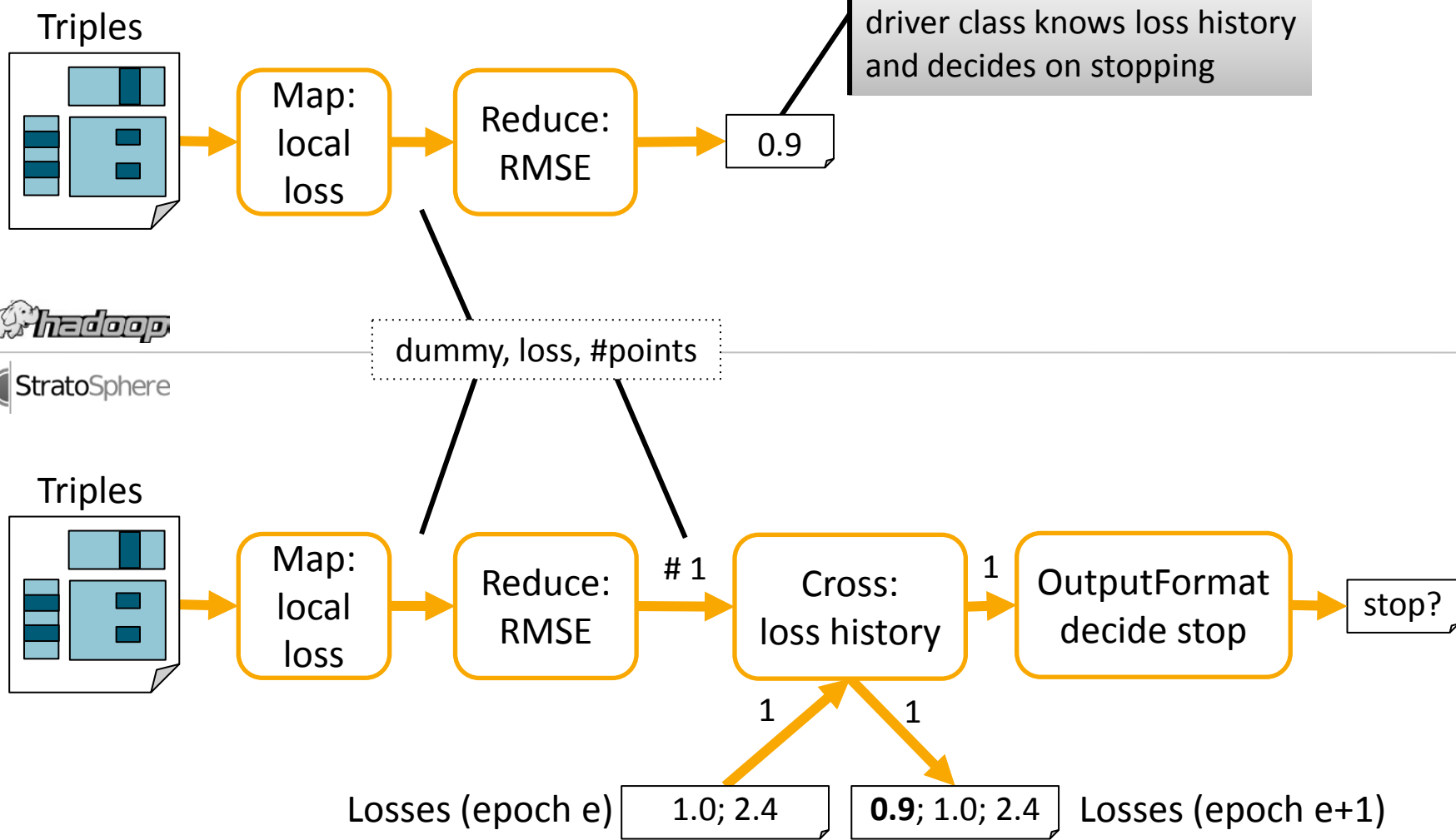


Stratosphere Plans

Calculate Loss (Training)

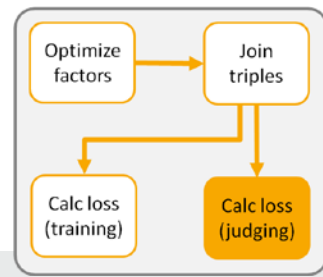


10

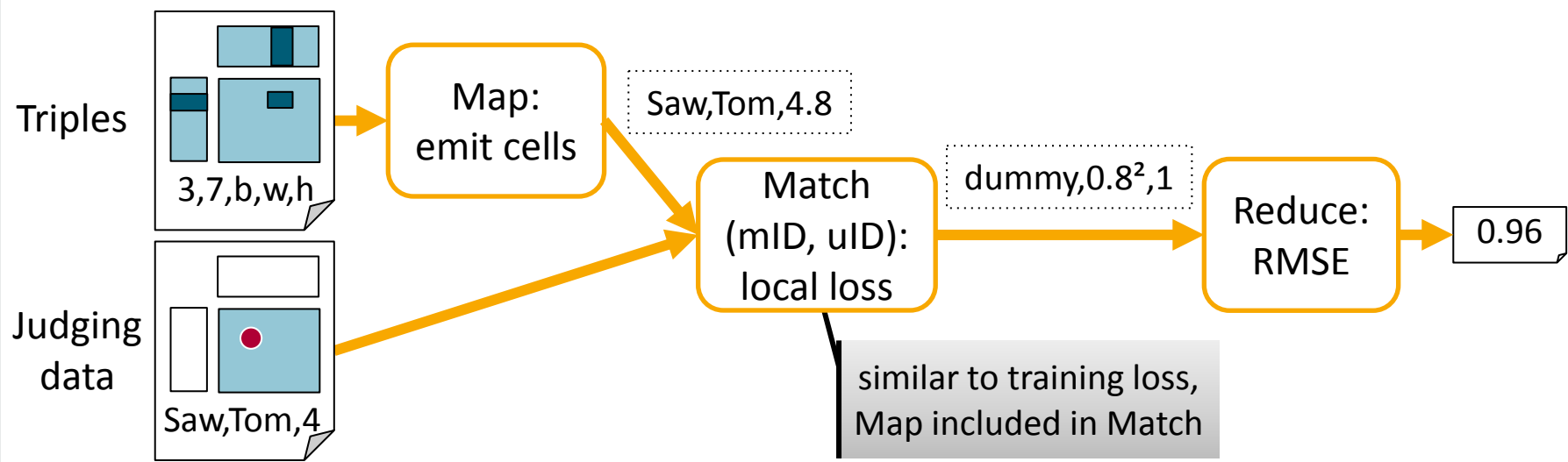
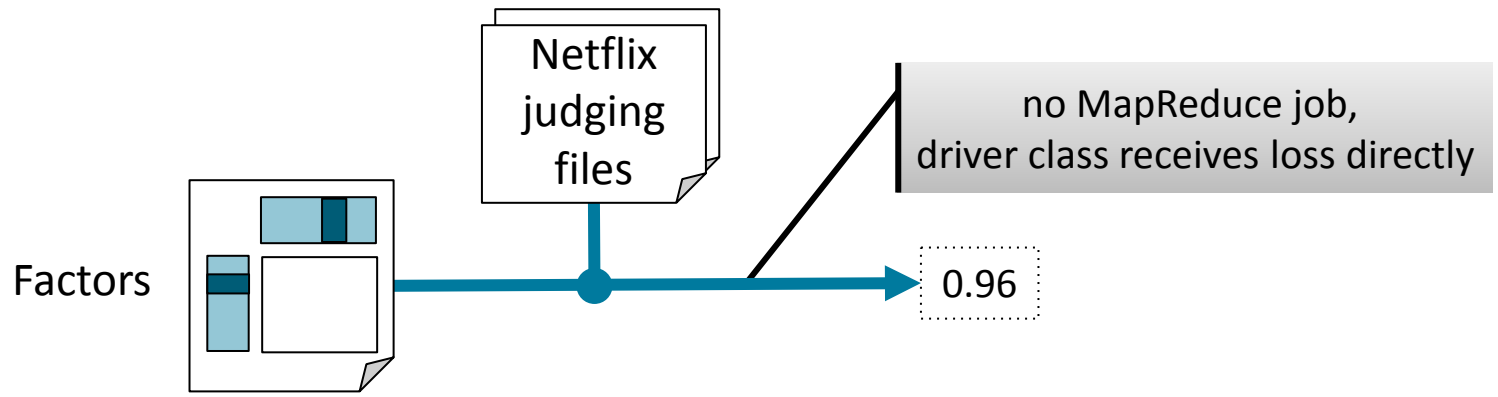


Stratosphere Plans

Calculate Loss (Judging)



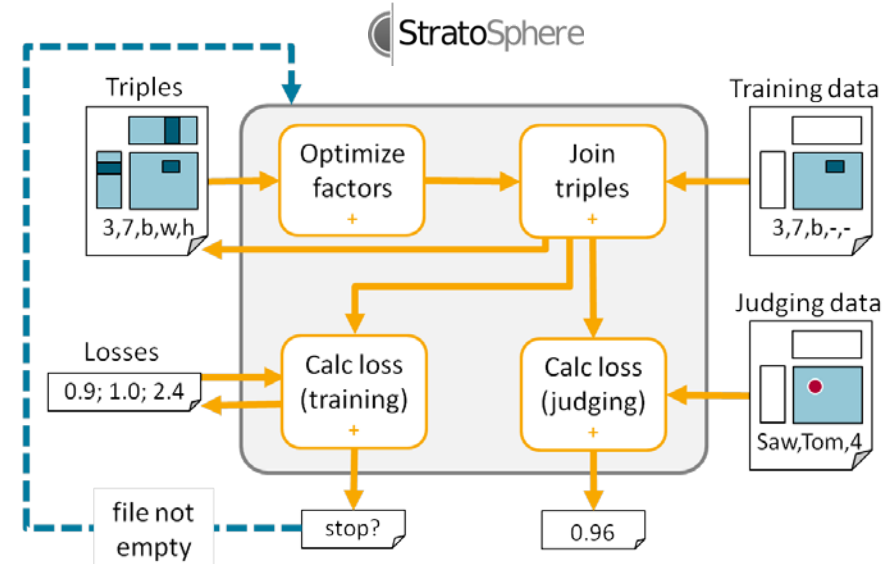
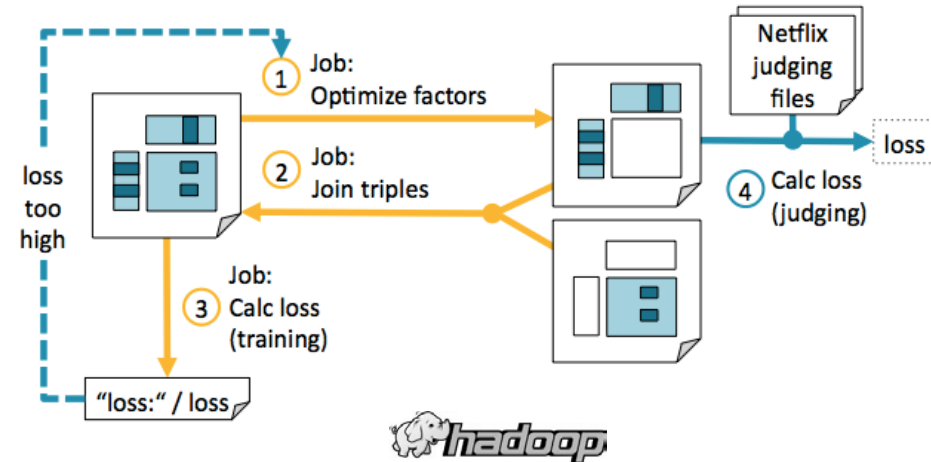
11



Comparison Jobs vs. Plans

12

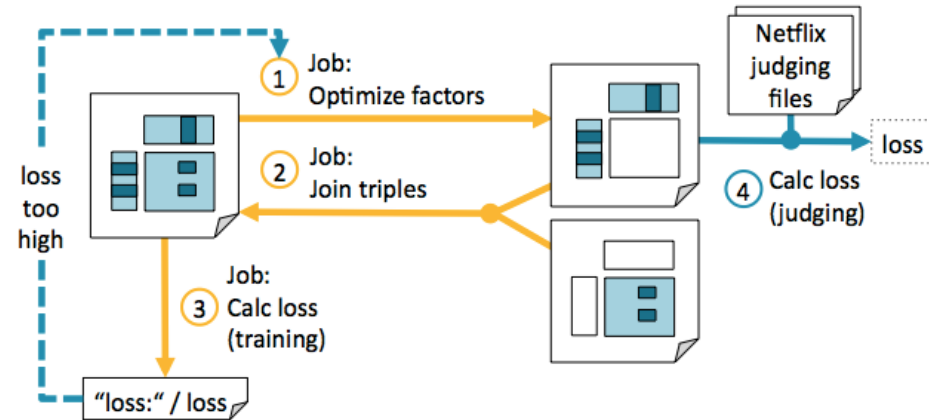
- Equal results without random...
 - Starting points
 - Training sequence
- Files
 - Factors not materialized
 - Separate factor files possible
 - Create new file for each iteration
- No efficient serialization between plans (yet)
 - Either parsing text file
 - Or use sequence file single-threaded



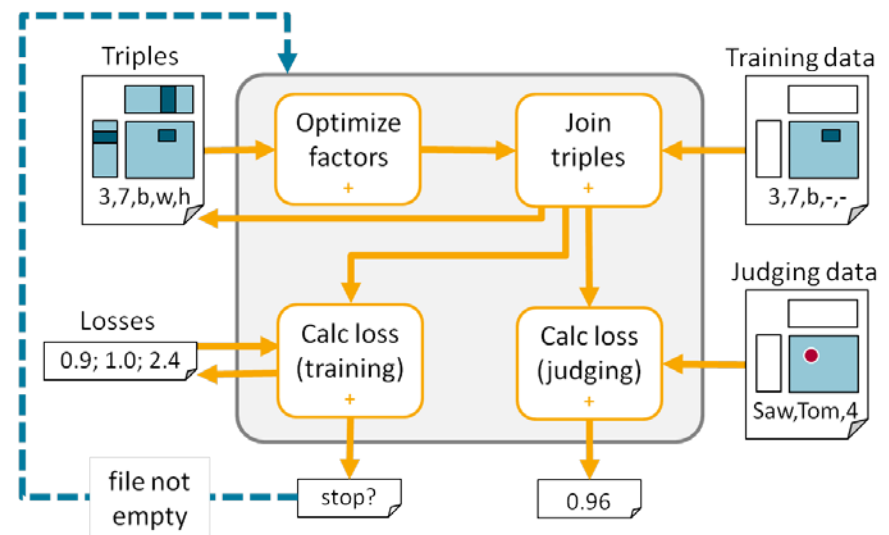
Comparison Data Schema

13

- 3#7 / TripleStorage
(Storage class is tagged union)
- Dummy key / LossStorage
- ➔ Getter, setter, toString()



- 3,7,b,w,h
- Dummy key, loss, #points
- ➔ Remember key places
- ➔ Reuse parts with configurations for different keys
- ➔ Composite keys possible



Comparison Preprocessing

14

Requires:

- Line format according to parameters
- Copy to HDFS
- Serialize factors and blocks



- Use Map file to write serialized values to HDFS

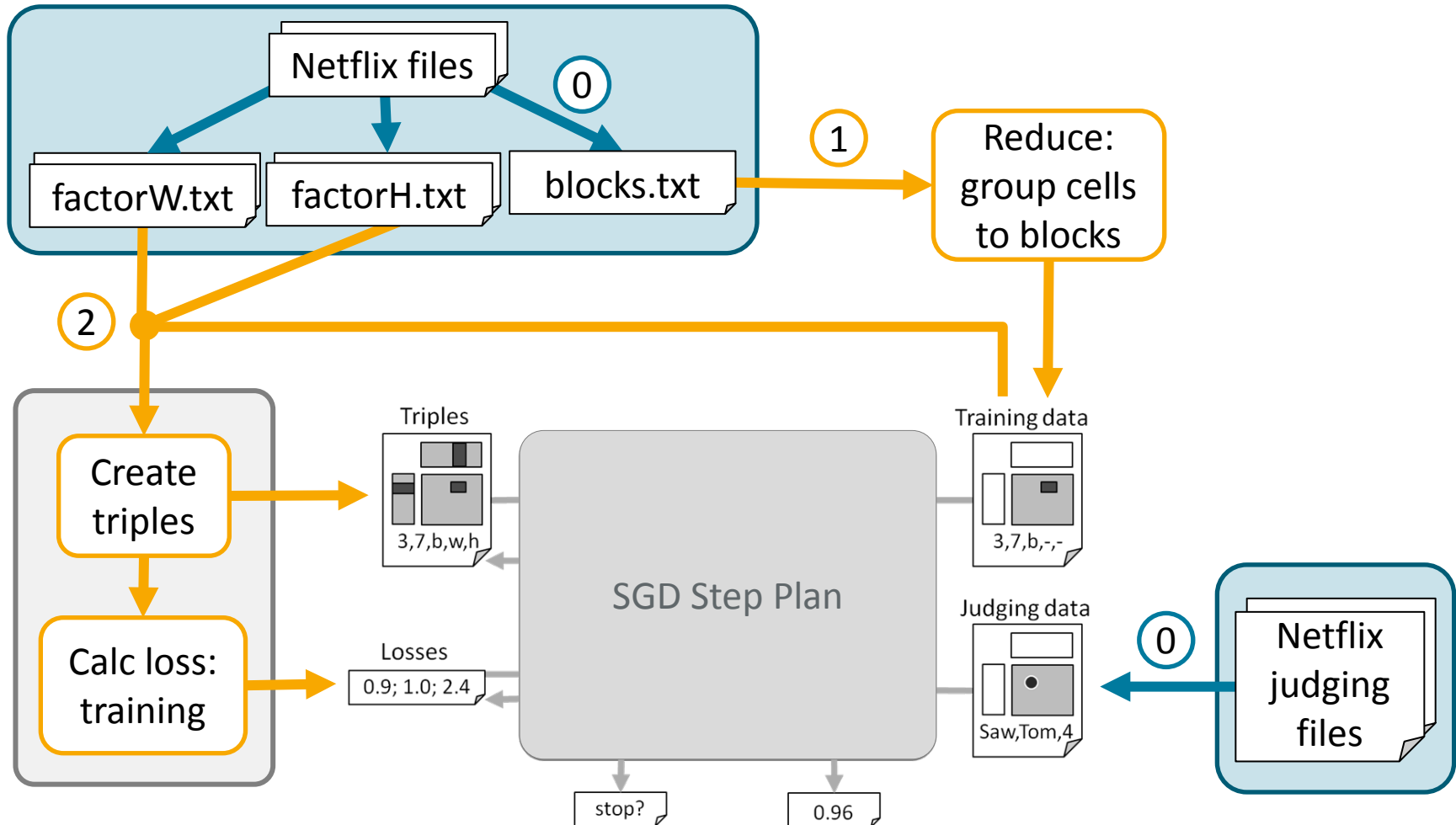


- Java process to define lines
- Shell script to move to HDFS
- Extra pacts to parse lines

Stratosphere Preprocessing

Define Line Format

15



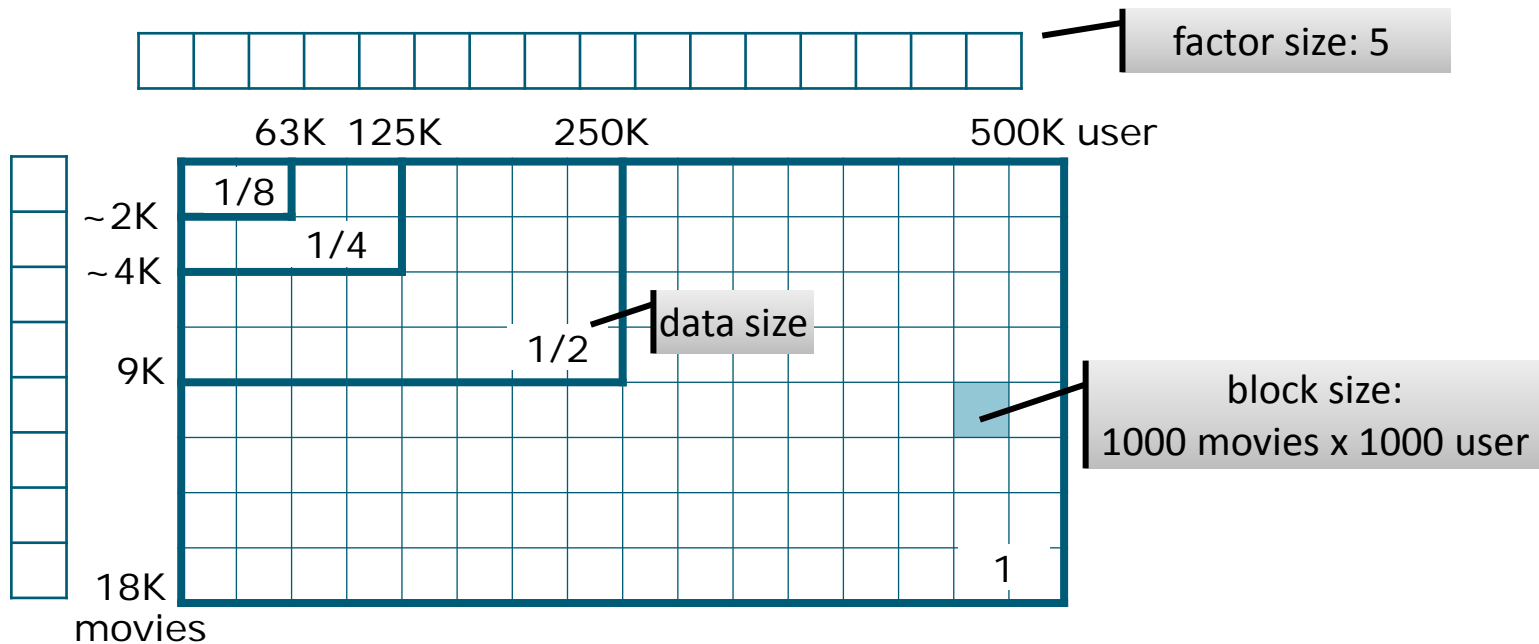
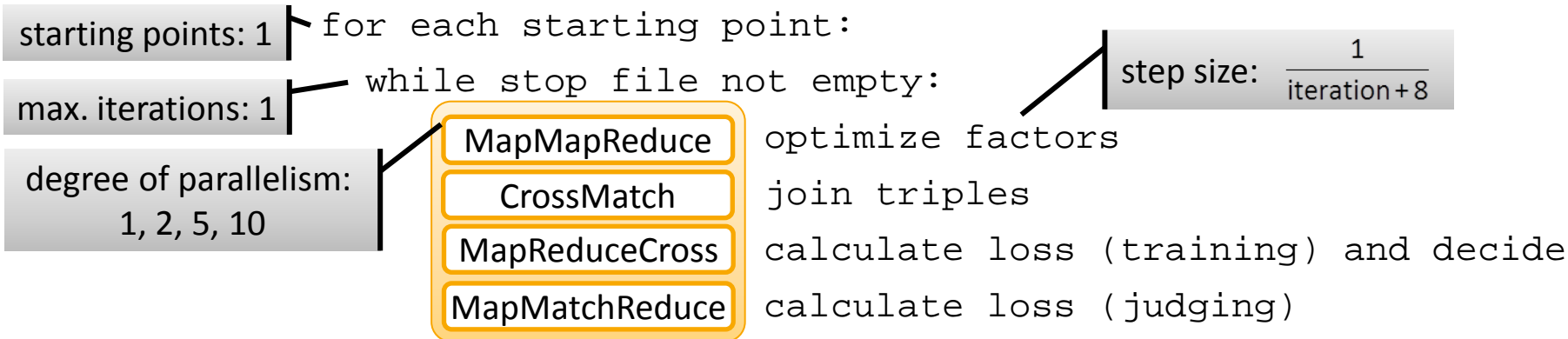
Suggestions

16

- Global aggregation of loss with dummy key:
 - Reducer with no key
 - Reducer with compiler hint `nrOfKeys = 1`
 - No sorting needed
- Provide `toString()` in `PactRecord`
- Provide getter, setter in `PactRecord` to encapsulate field numbers and class types
 - Keep configuration options (e.g. `reduce: average W or H`)
- Sequence files for sinks and sources
- Move log files to master node

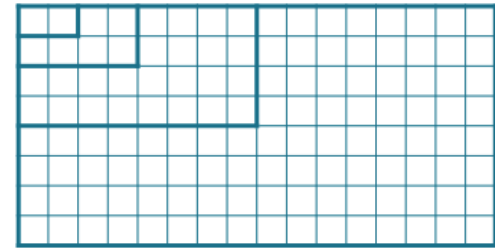
Evaluation Parameters for Netflix Data

17

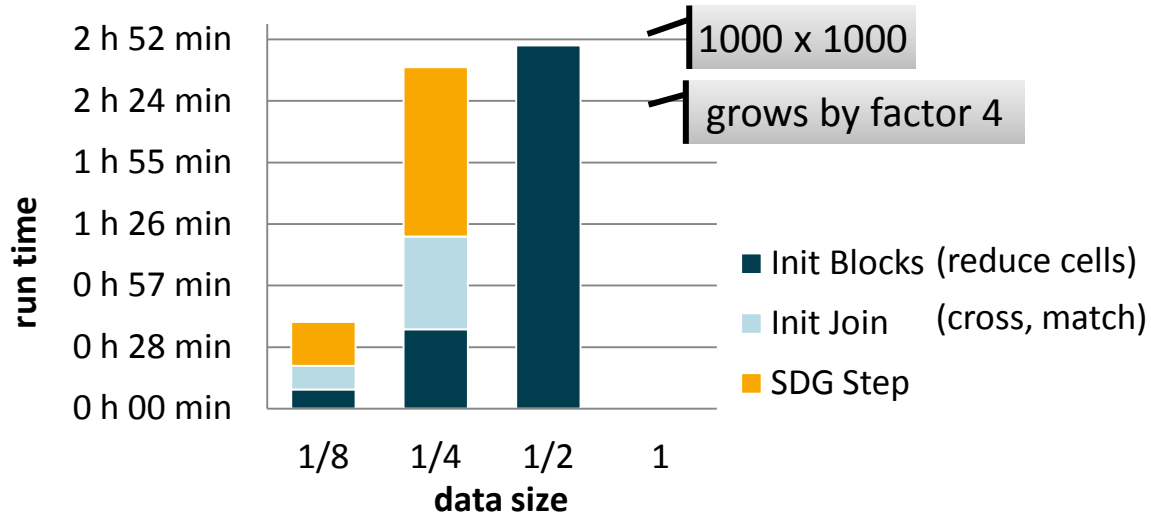
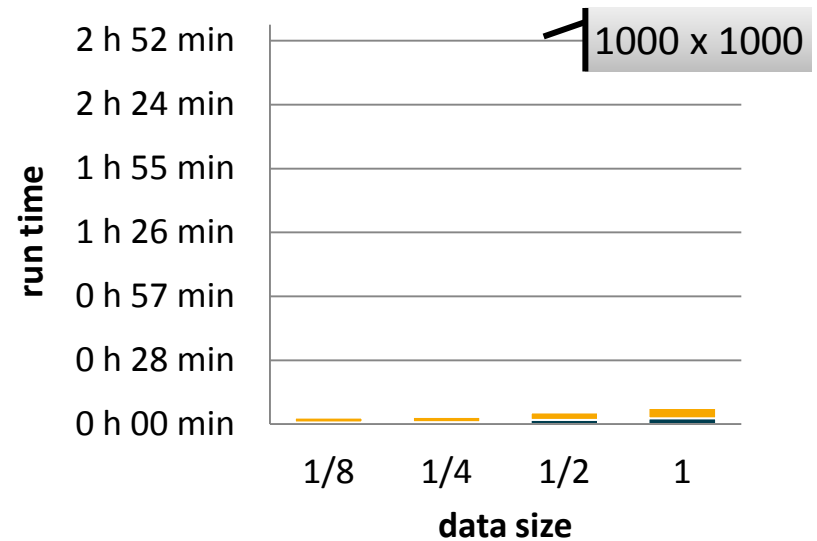
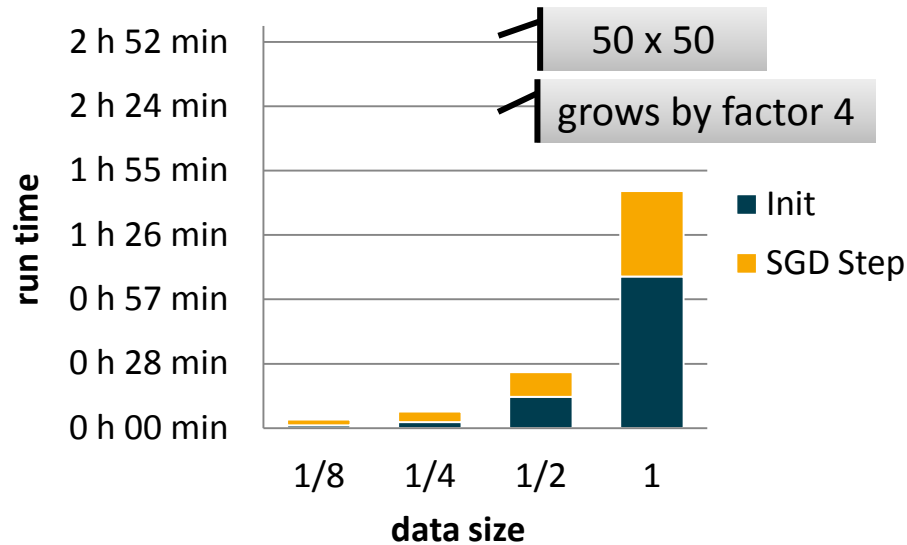


Evaluation

Run Time for Variable Data Size

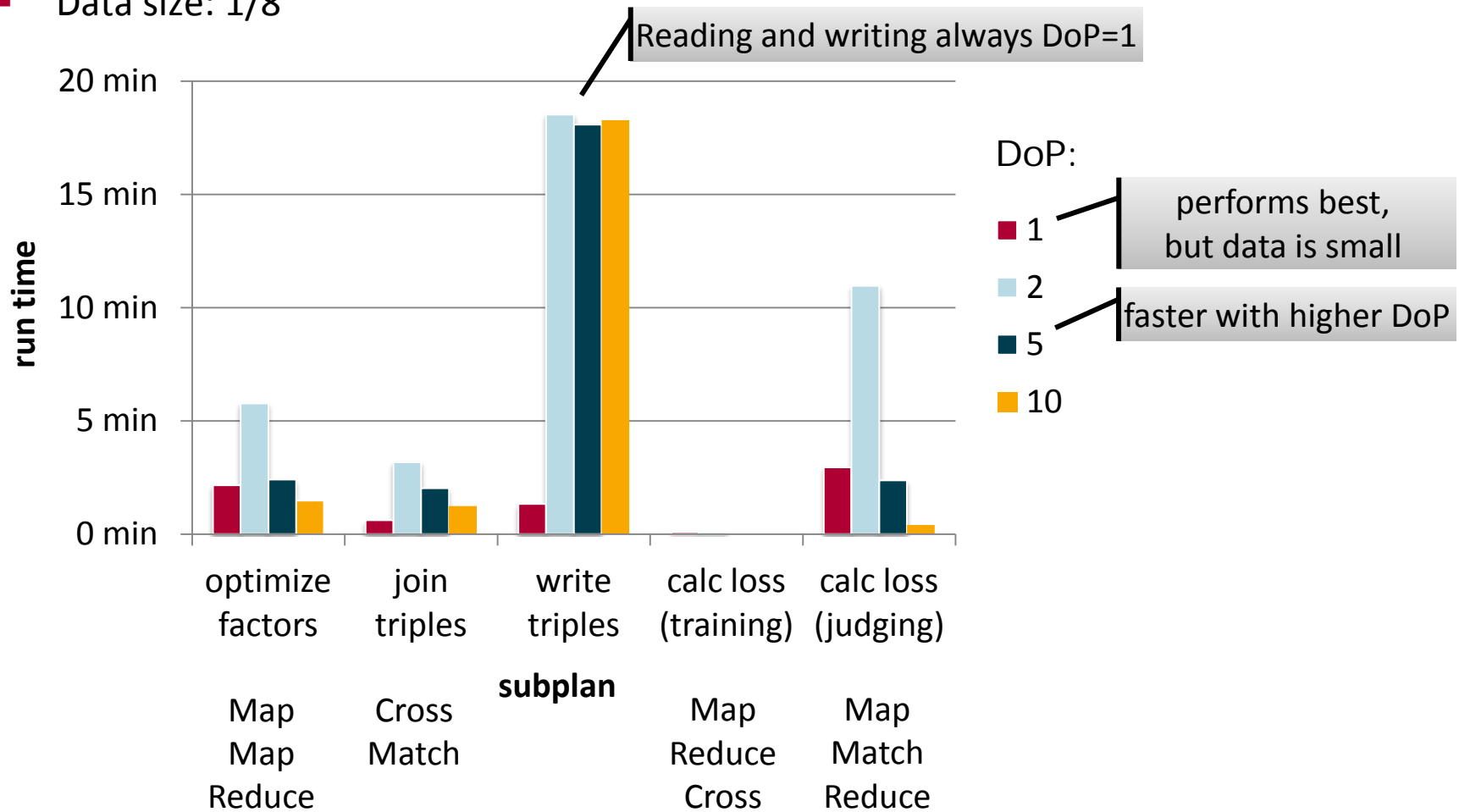


18



- 10 nodes vs. DoP = 8
- Usually 10 - 30 iterations
- No sequence file between plans

- Data size: 1/8



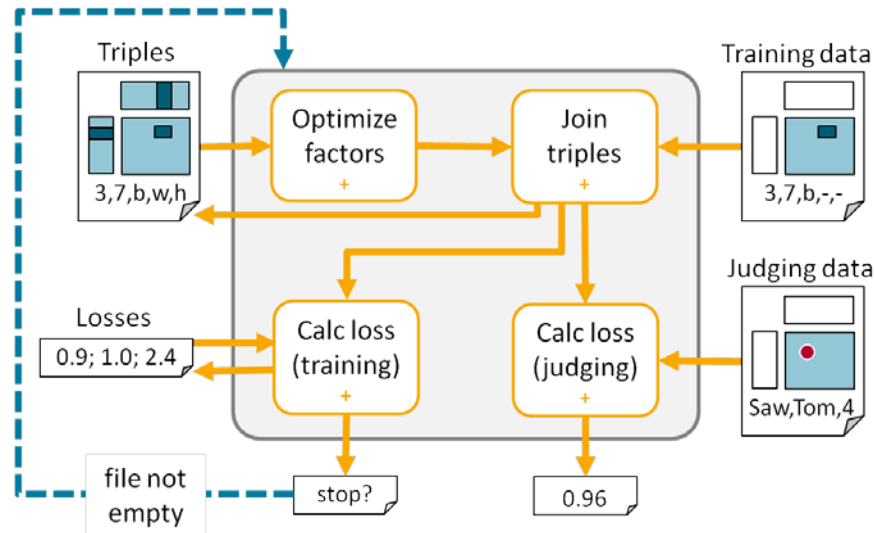
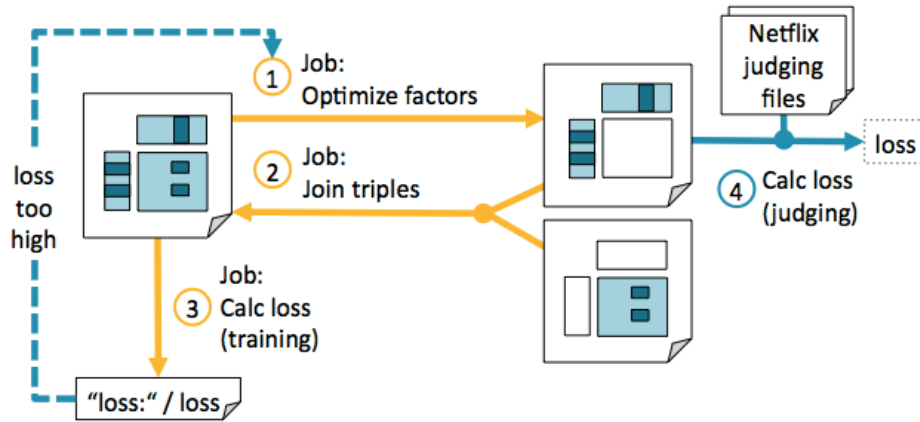
Outlook

20

- Join triples with Cross-Match vs. Match-Match
- Degree of parallelism: 10
- Inspect judging outcome: RMSE should be equal
- Evaluate DoP with bigger data

Summary

21



References

22

- Anand Rajaraman and Jeff Ullman. *Mining of Massive Datasets*. Cambridge University Press, 2010.
 - Rainer Gemulla, Peter J. Haas, Erik Nijkamp, and Yannis Sismanis. *Large-Scale Matrix Factorization with Distributed Stochastic Gradient Descent*. IBM Research Report RJ10481, March 2011.
- [1] <http://www.mathworks.de/matlabcentral/forums/27631/1/fff.png>, November 2011