

# Map Reduce / Workflow Systems II

Alexander Schönhuth



University Bielefeld  
May 28, 2020

# LEARNING GOALS TODAY

- ▶ Get to know idea of workflow systems and some examples
- ▶ Understand the definition of *communication cost*
- ▶ Understand the definition of *wall clock time*
- ▶ Get to know theory and intuition of *complexity theory* for MapReduce

# *Workflow Systems*

# WORKFLOW SYSTEMS: INTRODUCTION

- ▶ Workflow systems generalize MapReduce
- ▶ Just as much as MapReduce:
  - ▶ They're built on distributed file systems
  - ▶ They orchestrate large numbers of tasks with only small input provided by the user
  - ▶ They automatically handle failures
- ▶ In addition:
  - ▶ Single tasks can do other things than just Map or Reduce
  - ▶ Tasks interact in more complex ways

# WORKFLOW SYSTEMS: FLOW GRAPH

- ▶ A *function* represents arbitrary functionality within a workflow
  - ▶ and not just 'Map' or 'Reduce'
- ▶ Functions are represented as *nodes* of the *flow graph*
- ▶ Arcs  $a \rightarrow b$  for two functions  $a, b$  mean that the output of function  $a$  is provided to function  $b$  as input
- ▶ *Note:* The same function could be used by many tasks

# WORKFLOW SYSTEMS

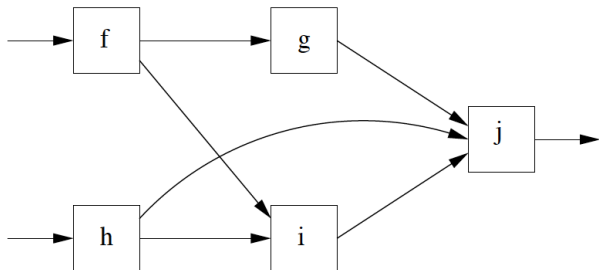


Figure: More complex workflow than MapReduce

Adopted from [mmds.org](http://mmds.org)

# WORKFLOW SYSTEMS: ACYCLIC FLOW GRAPH

- ▶ It is easier to deal with *acyclic flow graphs*
  - ▶ This means that one cannot return to functions
- ▶ *Blocking Property*: tasks only generate output upon completion
  - ▶ Blocking property easily applicable only in acyclic workflows
- ▶ *Simple Example of Workflow*: Cascades of Map-Reduce jobs

# POPULAR WORKFLOW SYSTEMS

- ▶ *Spark*: developed by UC Berkeley
- ▶ *TensorFlow*: Google's system, primarily developed for neural network computations
- ▶ *Pregel*: also by Google, for handling *recursive* (i.e. cyclic) workflows
- ▶ *Snakemake*: easy-to-use workflow system, inspired by MakeFile logic/functionality



# SPARK

- ▶ State-of-the-art workflow system:
  - ▶ Very efficient with failures
  - ▶ Very efficient in grouping tasks among nodes
  - ▶ Very efficient in scheduling execution of functions
- ▶ Basic concept: *Resilient Distributed Dataset (RDD)*
  - ▶ Generalizes key-value pair type of data: RDD is a file of objects of one type
  - ▶ *Distributed*: broken into chunks held at different nodes
  - ▶ *Resilient*: recoverable from losses of (even all) chunks
- ▶ *Transformations* (steps of functions) turn RDD into others
- ▶ *Actions* turn other data (from surrounding file system) into RDD's and vice versa

# SPARK: TRANSFORMATIONS

- ▶ *Map* takes a function as parameter and applies it to every element of an RDD, generating a new one
  - ▶ Turns one object into exactly another object, but not several ones
  - ▶ Remember: Map from MapReduce generates several key-value pairs from one object
- ▶ *Flatmap* is like Map from MapReduce, and generalizes it from key-value pairs to general object types
- ▶ *Filter* takes a predicate as input
  - ▶ Predicate is true or false for elements of RDD
  - ▶ So RDD is filtered for objects for which predicate applies
  - ▶ Yields a 'filtered RDD'

# SPARK: REDUCE AND RELATIONAL DATABASE OPERATIONS

- ▶ *Reduce* is an action, and takes as parameter a function that
  - ▶ applies to two elements of a particular type  $T$
  - ▶ returns one element of type  $T$
  - ▶ and is applied repeatedly until a single element remains
  - ▶ Works for associative and commutative operations
- ▶ Many *Relational Database Operations* are implemented in Spark:
  - ▶ Process RDD's reflecting tuples of relations
  - ▶ *Examples:* Join, GroupByKey

# SPARK: IMPLEMENTATION DETAILS

- ▶ Spark is similar like MapReduce in handling data (chunks are called *splits*)
- ▶ *Lazy evaluation* allows to apply several transformations consecutively to splits:
  - ▶ No intermediate formation of entire RDD's
  - ▶ Contradicts blocking property, because partial output is passed on to new functions
- ▶ *Resilience* (despite lazy evaluation) is maintained by *lineages of RDD's*
- ▶ Beneficial trade-off of more complex recovery of failures versus greater speed overall
  - ▶ Note that greater speed reduces probability of failures

# TENSORFLOW

- ▶ Open-source system developed (initially) by Google for machine-learning applications
- ▶ Programming interface for writing sequences of steps
- ▶ Data are *tensors*, which are multidimensional matrices
- ▶ Power comes from built-in operations applicable to tensors

# RECURSIVE WORKFLOWS

*Examples:*

- ▶ Calculating fixed-points ( $Mv = v$  for a matrix  $M$  and  $v$ ) by iterative application of  $M$  to  $v$
- ▶ Gradient descent, e.g. required in TensorFlow for determining optimal sets of parameters for machine learning models
- ▶ *Lack of blocking property:*
  - ▶ Flow graphs have cycles
  - ▶ Tasks may provide their output as input to other tasks whose output in turn results in more input to the first task
  - ▶ So generation of output only when task is done does not work
  - ▶ *Recovery from failures need to be reorganized*

# RECURSIVE WORKFLOWS: EXAMPLE

- ▶ Directed graph stored as relation  $E(X, Y)$ , listing arcs from  $X$  to  $Y$
- ▶ Want to compute relation  $P(X, Y)$ , listing paths from  $X$  to  $Y$
- ▶  $P$  is transitive closure of  $E$  (see below)
- ▶ *Algorithm:*
  - ▶ *Start:*  $P(X, Y) = E(X, Y)$
  - ▶ *Iteration:* Add to  $P$  tuples

$$\pi_{X,Y}(P(X, Z) \bowtie P(Z, Y)) \quad (1)$$

as pairs of nodes  $X$  and  $Y$  s.t. for some node  $Z$  there is path from  $X$  to  $Z$  and from  $Z$  to  $Y$

# TRANSITIVE CLOSURE: DEFINITION

DEFINITION [TRANSITIVE CLOSURE]:

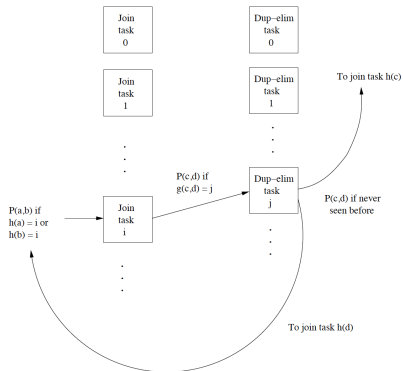
Let  $R(X, Y)$  be a relation.

- ▶  $R(X, Y)$  is *transitive* if  $(x, z) \in R$  and  $(z, y) \in R$  imply that  $(x, y) \in R$  as well
- ▶ The *transitive closure*  $\overline{R(X, Y)}$  of  $R(X, Y)$  is the *smallest set of tuples to be added to  $R(X, Y)$  that renders the resulting set of tuples transitive*



# EXAMPLE: TRANSITIVE CLOSURE

- ▶  $n$  Join tasks, corresponding to buckets of hash function  $h$
- ▶ Tuple  $P(a, b)$  is assigned to Join tasks  $h(a)$  and  $h(b)$
- ▶  $i$ -th Join tasks receives  $P(a, b)$ 
  - ▶ Store  $P(a, b)$  locally
  - ▶ If  $h(a) = i$  look for tuples  $P(x, a)$  and produce  $P(x, b)$
  - ▶ If  $h(b) = i$  look for tuples  $P(b, y)$  and produce  $P(a, y)$

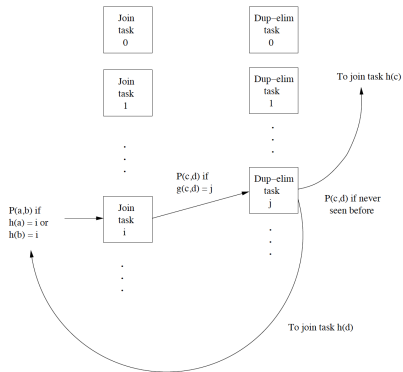


Transitive closure by recursive tasks

Adopted from [mmds.org](http://mmds.org)

# RECURSIVE WORKFLOWS: EXAMPLE

- ▶  $m$  Dup-elim tasks, corresponding to buckets of hash function  $g$
- ▶  $P(c, d)$  (as output of Join task) is sent to Dup-elim task  $j = g(c, d)$
- ▶ Dup-elim task  $j$  checks whether  $P(c, d)$  was received before
  - ▶ If *yes*,  $P(c, d)$  is ignored (and not stored)
  - ▶ If *not*,  $P(c, d)$  is stored locally,
  - ▶ *and* sent to Join tasks  $h(c)$  and  $h(d)$

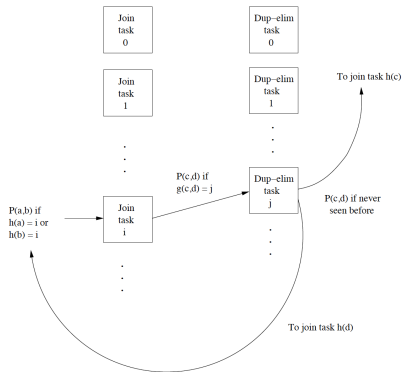


Transitive closure by recursive tasks

Adopted from [mmds.org](http://mmds.org)

# RECURSIVE WORKFLOWS: EXAMPLE

- ▶ Every Join task has  $m$  output files
- ▶ Every Dup-elim task has  $n$  output files
- ▶ Initially, tuples  $E(a, b)$  are sent to Dup-elim tasks  $g(a, b)$



Transitive closure by recursive tasks

Adopted from [mmds.org](http://mmds.org)

# RECURSIVE WORKFLOWS: FAILURE HANDLING

- ▶ *Iterated MapReduce*: Application is repeated execution / sequence of MapReduce job(s) (“HaLoop”)
- ▶ *Spark Approach*: Lazy evaluation, lineage mechanisms, option to store intermediate results
- ▶ *Bulk Synchronous Systems*: Graph-based model using “periodic checkpointing”

# BULK SYNCHRONOUS SYSTEMS: PREGEL

- ▶ System views data as *graph*:
  - ▶ *Nodes* (roughly) reflect tasks
  - ▶ *Arcs*: from nodes whose output (messages) are input to other nodes
- ▶ *Supersteps*:
  - ▶ All messages received by any of the nodes from the previous superstep are processed
  - ▶ All messages generated are sent to their destinations
- ▶ *Advantage*: Sending messages means communication costs, bundling them reduces costs
- ▶ *Failure Management*: Checkpointing entire computation by making copy after each superstep
- ▶ May be beneficial to checkpoint periodically after number of supersteps

# SNAKEMAKE

- ▶ Create *reproducible* and *scalable* data analyses
- ▶ Workflows described in human readable, Python based language
- ▶ Seamlessly scale to server, cluster, grid and cloud environments
- ▶ Integrating descriptions of required software, deployable to any execution environment

# *The Communication-Cost Model*

# COMMUNICATION COST

## Situation

- ▶ Algorithm implemented by acyclic network of tasks:
  - ▶ Map tasks feeding Reduce tasks
  - ▶ Cascade of several MapReduce jobs
  - ▶ More general workflow structure (e.g. Fig. 1)

## DEFINITION [COMMUNICATION COST]:

- ▶ The *communication cost of a task* is the size of the input it receives
- ▶ The *communication cost of an algorithm* is the sum of the communication costs of its tasks



# COMMUNICATION COST

## Why Communication Cost?

- ▶ Computing communication cost is the way to measure the complexity of distributed algorithm
- ▶ Neglect time necessary for tasks to execute
- ▶ Importance of communication cost:
  - ▶ Tasks tend to be simple (often linear in size of input)
  - ▶ Interconnect speed of compute cluster (typically 1 Gbit/sec) slow compared with speed processors execute instructions
  - ▶ Often there is competition for the interconnect when several nodes are communicating
  - ▶ Moving data from disk to memory may exceed runtime

## Why not Output Size?

- ▶ Output often is input to another task anyway
- ▶ Output rarely large in comparison with input or intermediate data

# REMINDER: NATURAL JOIN

**Natural Join:**  $R(A, B) \bowtie S(B, C)$

- ▶ **Map:** For each tuple  $t = (a, b)$  from  $R$ , generate key-value pair  $(b, (R, a))$ . For each tuple  $(b, c)$  from  $S$ , generate  $(b, (S, c))$ .
- ▶ **Reduce:** After grouping, each key value  $b$  has list of values being either of the form  $(R, a)$  or  $(S, c)$ 
  - ▶ Construct all pairs of values where first component is like  $(R, a)$  and second component is like  $(S, c)$ , yielding triples  $(b, (R, a), (S, c))$
  - ▶ Turn triples into triples  $(a, b, c)$  being output

# COMMUNICATION COST: NATURAL JOIN EXAMPLE

Suppose we are joining  $R(A, B) \bowtie S(B, C)$  with  $R, S$  of sizes  $r$  and  $s$ .

- ▶ *Map*: Chunks of files  $R, S$  are input to Map tasks
  - ☞ communication cost of Map is  $r + s$  (in practice mostly disk to memory)
- ▶ *Reduce*: Input to Reduce tasks is all  $(r + s)$  many key-value pairs generated by Map tasks
  - ☞ communication cost for Reduce is  $O(r + s)$
- ▶ *Output of Reduce* could be much larger than  $O(r + s)$  (up to  $O(rs)$ ), depending on how many tuples are to be generated for each key  $b$

# COMMUNICATION COST EXAMPLE: $R(A, B) \bowtie S(B, C)$

Let sizes of relations  $R$  and  $S$  be  $r$  and  $s$ .

## Map

- ▶ Each chunk of the files holding  $R$  and  $S$  is fed to one task
  - ↳ Communication cost is  $r + s$
- ▶ Nodes hold chunks already from file distribution step: no internode communication, only disk-to-memory costs
- ▶ All Map tasks perform a simple transformation, so only negligible computation cost
- ▶ Output about as large as input

# COMMUNICATION COST EXAMPLE: $R(A, B) \bowtie S(B, C)$

Let sizes of relations  $R$  and  $S$  be  $r$  and  $s$ .

## Reduce

- ▶ Receives and divides input into tuples from  $R$  and  $S$
- ▶ For each key, pairs each tuple from  $R$  with the ones from  $S$
- ▶ Output size can vary: can be larger or smaller than  $O(r + s)$ 
  - ▶ Many different B-values: output is small
  - ▶ Few B-values: output much larger
- ▶ Output large: computation cost could be much larger than  $O(r + s)$
- ▶ Often output is further subsequently aggregated at further nodes
  - ☞ Communication cost greater than computation cost

# WALL-CLOCK TIME

DEFINITION [WALL-CLOCK TIME]:

The *wall-clock time* is defined to be the time for the entire parallel algorithm to finish.

*Example:* Careless reasoning could make one assign all tasks to one node, which minimizes communication cost. But the wall-clock time is (likely to be) at its maximum.

## EXAMPLE: MULTIWAY JOIN

Consider computing  $R(A, B) \bowtie S(B, C) \bowtie T(C, D)$ . For simplicity, let  $p$  be the probability that *both* an  $R$ - and an  $S$ -tuple agree on  $B$  and that an  $S$ - and a  $T$ -tuple agree on  $C$

- ▶ Joining  $R$  and  $S$  first
  - ▶ Communication cost is  $O(r + s)$  (see before)
  - ▶ Size of output is  $prs$
  - ▶ Hence joining  $R \bowtie S$  with  $T$  is  $O((r + s) + (t + prs))$
- ▶ Joining  $S$  and  $T$  first analogously yields  $O((s + t) + (r + pst))$

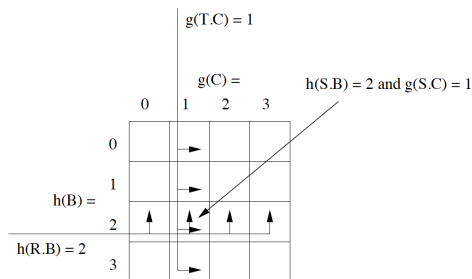
## $R(A, B) \bowtie S(B, C) \bowtie T(C, D)$ IN ONE MAPREDUCE

Let  $p$  be the probability that *both* an  $R$ - and an  $S$ -tuple agree on  $B$  and that an  $S$ - and a  $T$ -tuple agree on  $C$ .

- ▶ Hash  $B$ - and  $C$ -values, using functions  $h$  and  $g$
- ▶ Let  $b$  and  $c$  be the number of buckets for  $h$  and  $g$
- ▶ Let  $k$  be the number of Reducers; require that  $bc = k$ 
  - ▶ Each reducer corresponds to a pair of buckets
  - ▶ Reducer corresponding to bucket pair  $(i, j)$  joins tuples  $R(u, v), S(v, w), T(w, x)$  whenever  $h(v) = i, g(w) = j$
- ▶ Hence Map tasks send  $R$ - and  $T$ -tuples to more than one reducer
  - ▶  $R$ -tuples  $R(u, v)$  go to all reducers  $(h(v), y)$ 
    - ↳ goes to  $c$  reducers
  - ▶  $T$ -tuples  $T(w, x)$  go to all reducers  $(z, g(w))$ 
    - ↳ goes to  $b$  reducers



# MULTIWAY JOIN: ONE MAPREDUCE II



Sixteen reducers for a 3-way join

Adopted from [mmds.org](http://mmds.org)

- ▶  $h(v) = 2, g(w) = 1$
- ▶  $S$ -tuple  $S(v, w)$  goes to reducer for key  $(2, 1)$
- ▶  $R$ -tuple  $R(u, v)$  goes to reducers for keys  $(2, 0), \dots, (2, 3)$
- ▶  $T$ -tuple  $T(w, x)$  goes to reducers for keys  $(1, 0), \dots, (1, 3)$

# MULTIWAY JOIN: ONE MAPREDUCE III

## Communication cost:

- ▶ Moving tuples to proper reducers is sum of
  - ▶  $s$  to send tuples  $S(v, w)$  to  $(h(v), g(w))$
  - ▶  $rc$  to send tuples  $R(u, v)$  to  $(h(v), y)$  for each of the  $c$  possible  $g(w) = y$
  - ▶  $bt$  to send tuples  $T(w, x)$  to  $(z, g(w))$  for each of the  $b$  possible  $h(b) = z$
- ▶ Additional (constant) cost  $r + s + t$  to make each tuple input to one of the Map tasks (constant)

# MULTIWAY JOIN: ONE MAPREDUCE III

## Communication cost:

- ▶ *Goal:* Select  $b$  and  $c$ , subject to  $bc = k$ , to minimize  $s + cr + bt$
- ▶ Using Lagrangian multiplier  $\lambda$  yields to solve for
  - ▶  $r - \lambda b = 0$
  - ▶  $t - \lambda c = 0$
- ▶ It follows that  $rt = \lambda^2 bc$ , that is  $rt = \lambda^2 k$ , yielding further  $\lambda = \sqrt{\frac{rt}{k}}$
- ▶ So, minimum communication cost at  $c = \sqrt{\frac{kt}{r}}$  and  $b = \sqrt{\frac{kr}{t}}$
- ▶ Substituting into  $s + cr + bt$  yields  $s + 2\sqrt{krt}$
- ▶ Adding  $r + s + t$  yields  $r + 2s + t + 2\sqrt{krt}$ , which is usually dominated by  $2\sqrt{krt}$

# MATERIALS / OUTLOOK

- ▶ See *Mining of Massive Datasets*, chapter 2.4–2.6
- ▶ For deepening your understanding, voluntary *homework*: please read through 2.6.7
- ▶ As usual, see <http://www.mmds.org/> in general for further resources
- ▶ Next lecture: “Mining Data Streams”
  - ▶ See *Mining of Massive Datasets* 4.1–4.7