

Functional Models of Hadoop MapReduce with Application to Scan

Kiminori Matsuzaki
Kochi University of Technology



Brief History of (Hadoop) MapReduce

- 2004: Google proposed MapReduce [OSDI 2004]
- 2004-2006: Open-Source MapReduce in Nutch
- 2006-: Hadoop project
- 2011 Dec.: Hadoop 1.0.0
- 2012: Industry standard in distributed processing
- 2013 Oct.: Hadoop 2.2.0 (first stable ver. 2.x)

[<http://research.yahoo.com/files/cutting.pdf>]

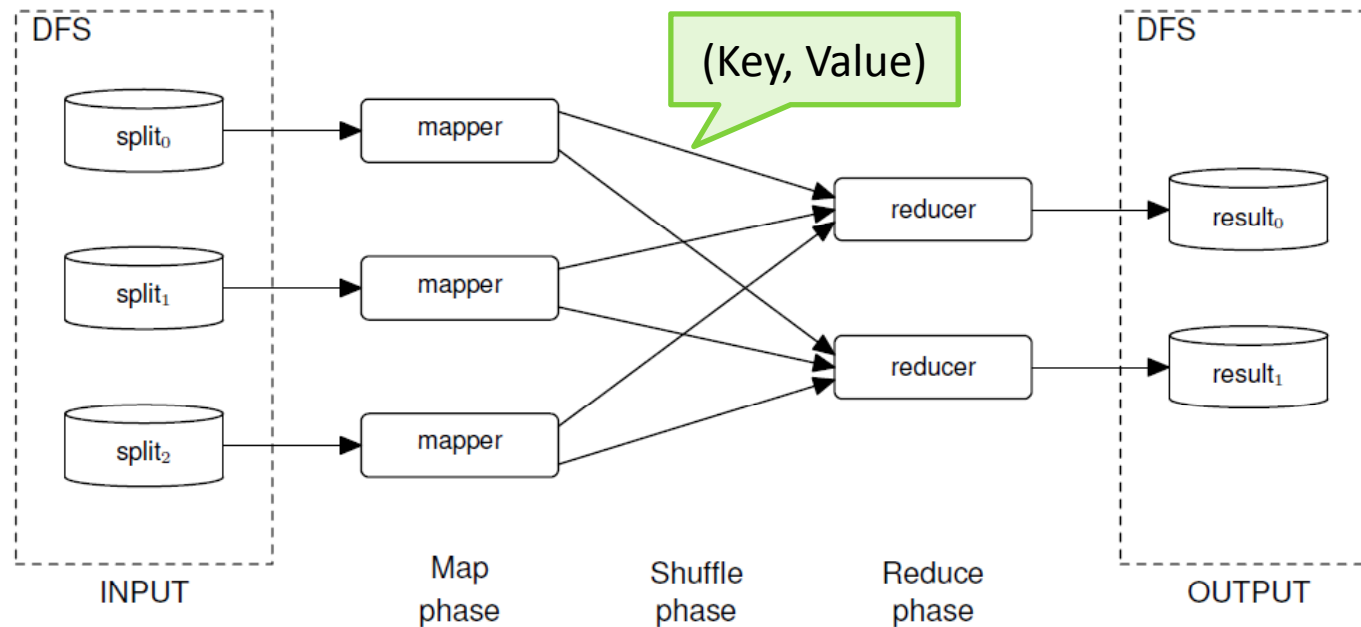
[<http://www.guruzon.com/6/introduction/map-reduce/history-of-map-reduce/>]

[<http://hadoop.apache.org/releases.html>]



MapReduce in a NutShell

- 3 phases, 2 user-defined functions





Misunderstandings

Functional
Community

DB
Community

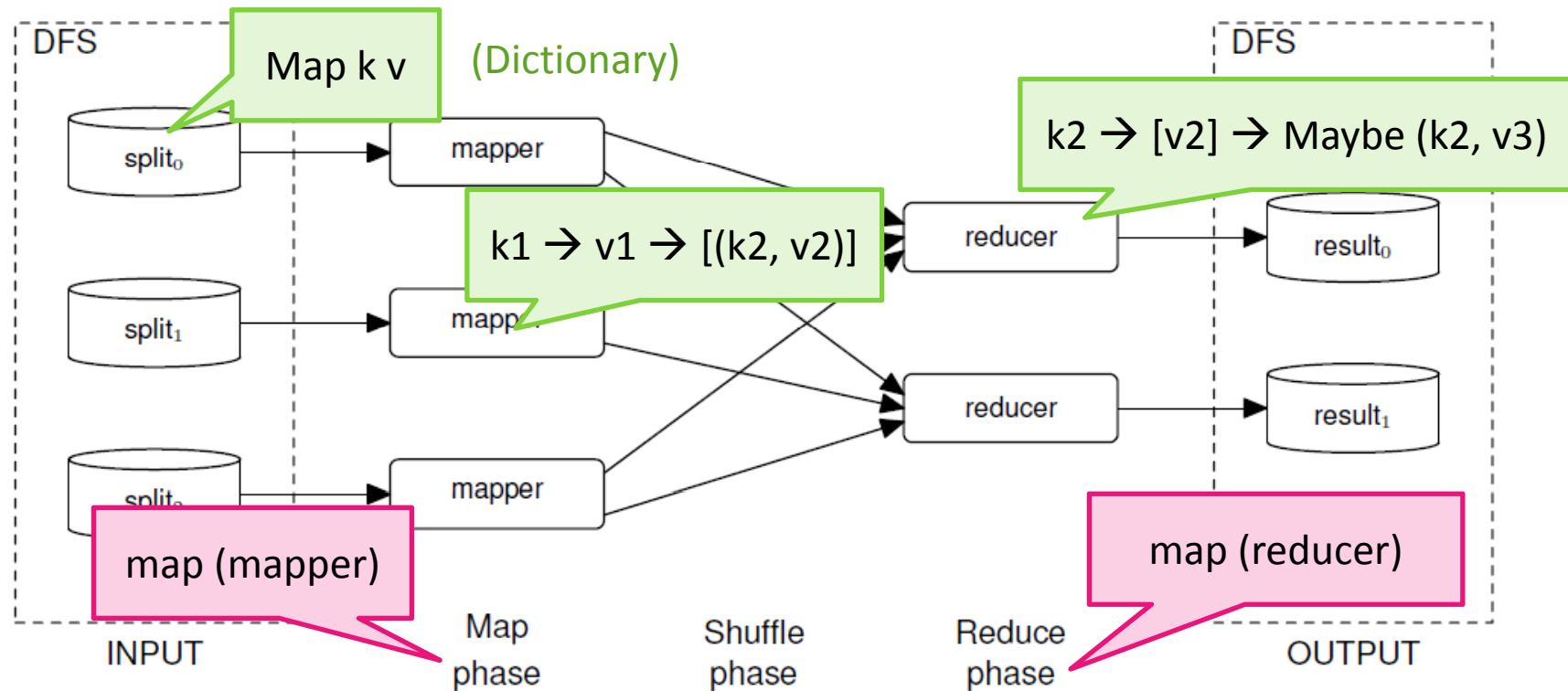
- The OSDI paper said: “map/reduce were inspired from those in functional programming”
- Map/reduce in MapReduce differs in terms of
 - target data (set-like data, not lists or trees)
 - how they work (map/reduce are applied independently)
 - no associativity needed in reduce
- “MapReduce: A major step backwards” (2008)
 - No indexing, poor impl., DBMS-incompatibility, etc.

A functional model describes **clearly** the computation of the framework, especially by using the **types**

- R. Lämmel:
 - “Google’s MapReduce programming model -- revisited.”
 - Science of Computer Programming*, 2008
 - Provides a functional model of Google’s MapReduce
 - Model is written in Haskell



Lämmel's Functional Model





Why Functional Model Matters?

- Understanding the computation
 - Avoid misunderstandings
- Proof of Correctness
 - Developing functional code to check
 - Proof using Coq (Related work [Ono 2011], [Jiang 2014])
- Program Calculation
 - Developing program-transformation rules
- Cost Model
 - Performance tuning (Related work [Dörre 2014])

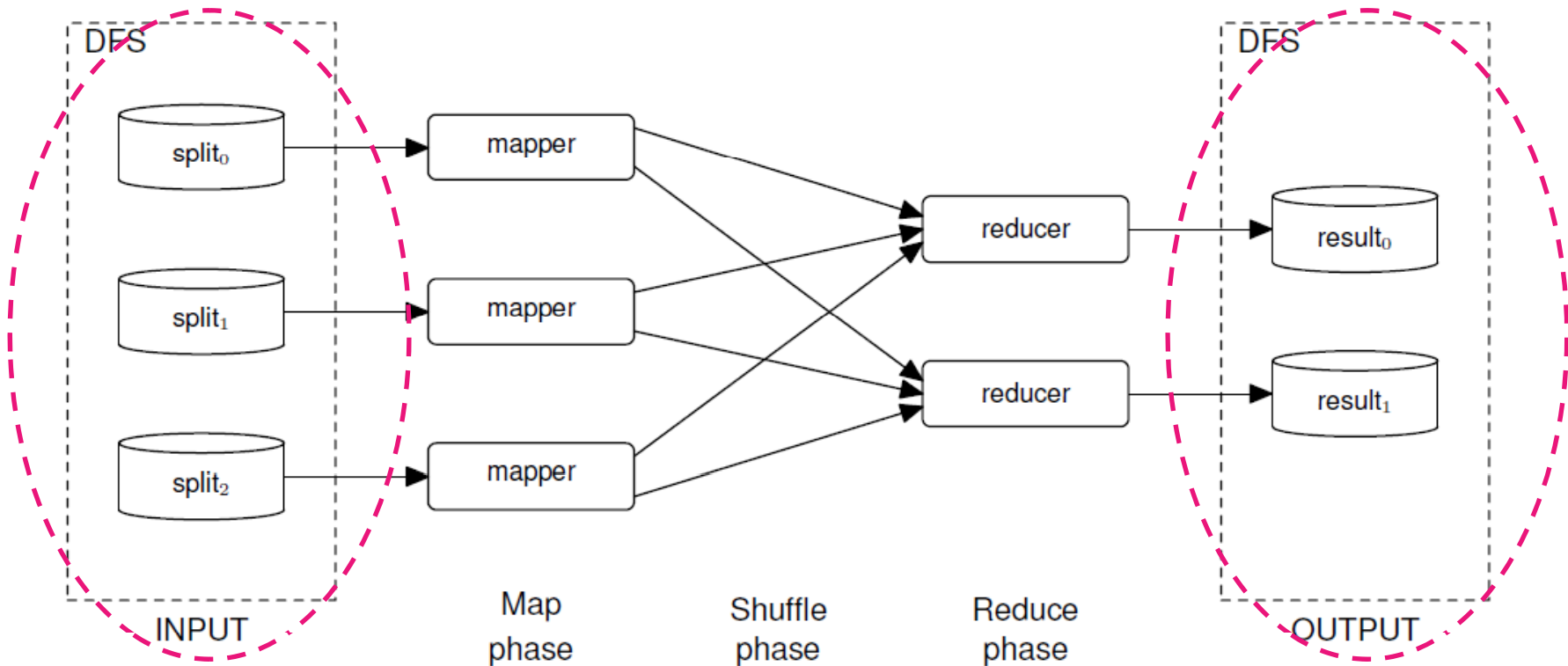


Contributions In The Paper

- Two functional models of Hadoop MapReduce
 - Low-level model (based on implementation)
 - Nested input/output
 - Stateful mapper/reducer
 - Detailed modeling of Shuffling phase
 - High-level model (user-friendly specification)
 - for “secondary-sorting” technique
- Scan (prefix-sums) algorithm on the models
 - Three-phase algorithm (L-reduce, G-scan, L-scan)
 - BSP-based 2-superstep algorithm



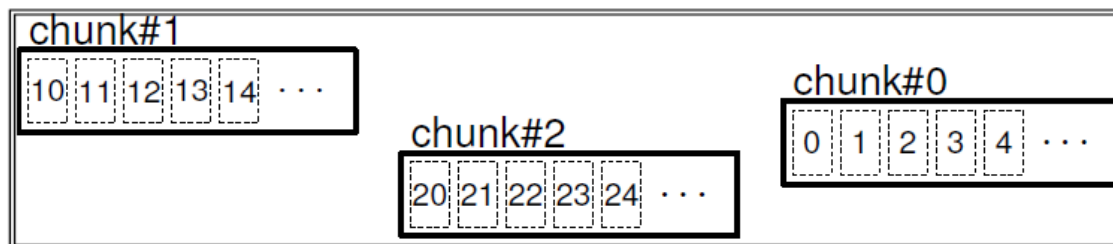
Nested Input/Output in Hadoop





Nested Input/Output in Hadoop

- Data → Split → Record

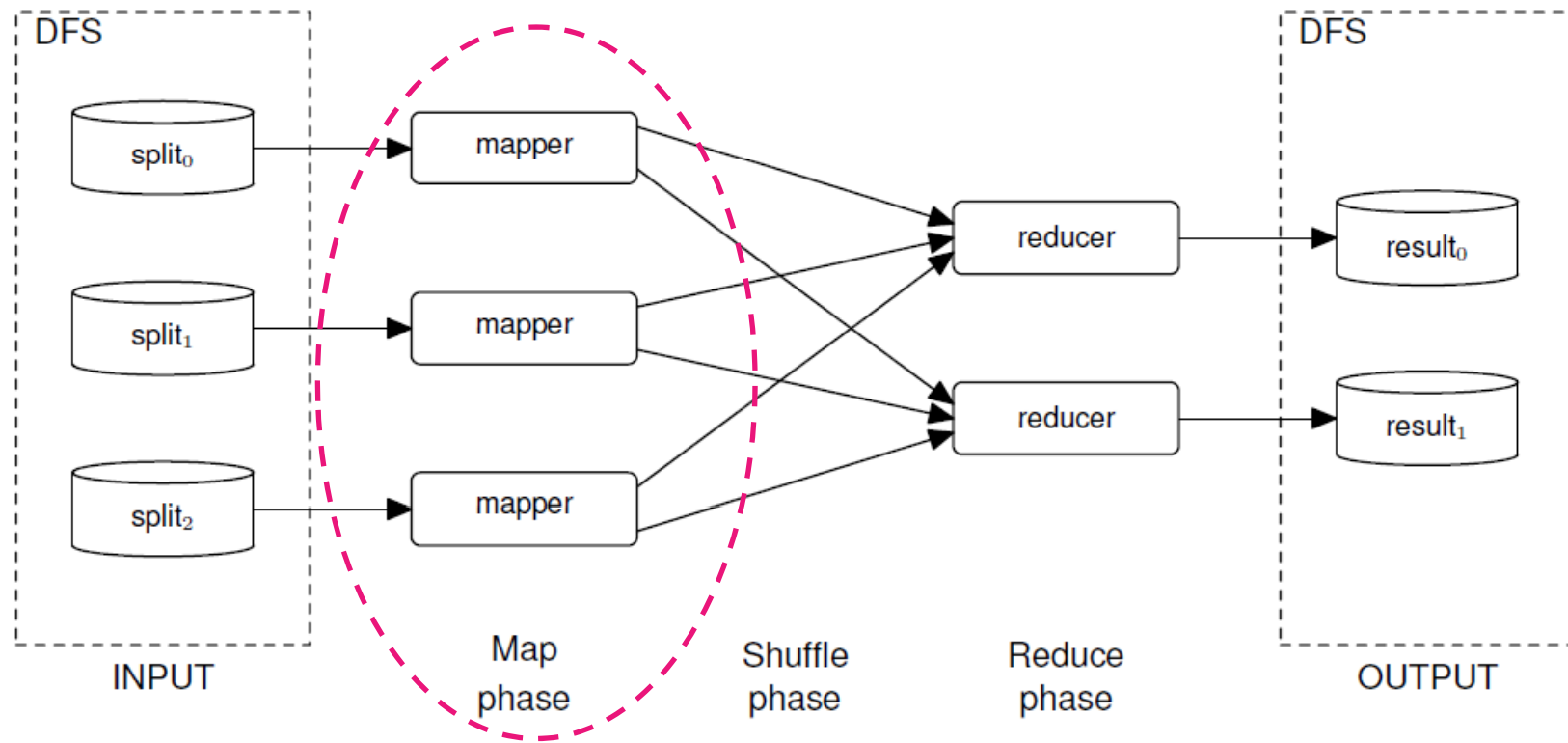


▭ The whole DFS ▭ Splits (unordered) ▭ Records (ordered)

- Split → Unordered → Bag (multi-set)
 - Parallel (independently / order may change)
- Record → Ordered → List
 - Sequential (one-by-one / order preserved)



Mapper Class in Hadoop





Mapper Class in Hadoop

- Definition in Hadoop

```
class Mapper {
  void setup(Context);
  void map(Key, Value, Context);
  void cleanup(Context);

  void run(Context c) {
    setup(c);
    for (kv : split) {
      map(kv.k, kv.v, c);
    }
    cleanup(c);
  }
}
```

- A simple case == map

```
mkMapper1 :: ((k1, v1) → [(k2, v2)])           -- f_map
              → [(k1, v1)] → [(k2, v2)]       -- input/output
mkMapper1 f_map = flatten ∘ map f_map
```

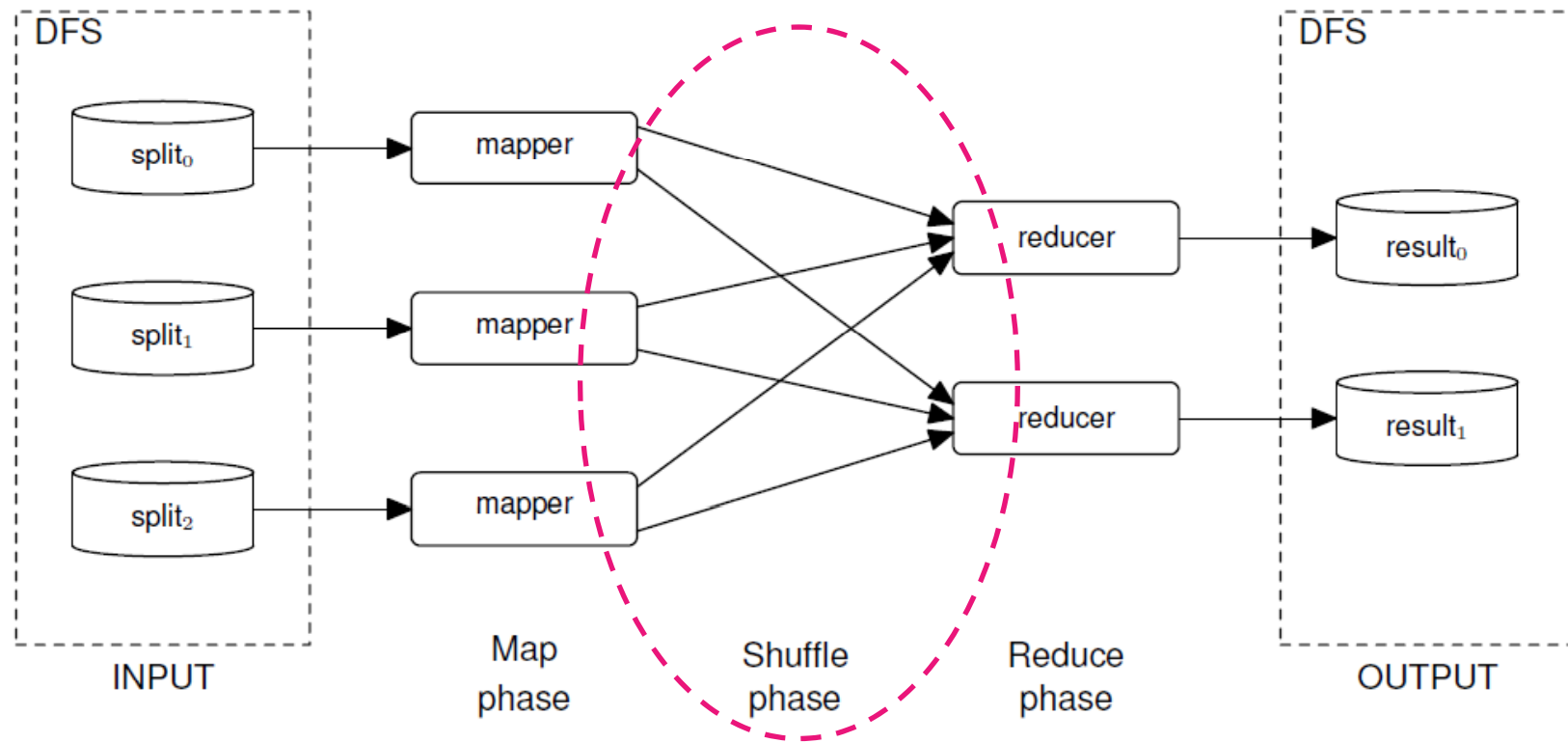
- A stateful case == foldl

```
mkMapper2 :: att                               -- f_setup
              → (att → (k1, v1) → att)        -- f_map
              → (att → [(k2, v2)])            -- f_cleanup
              → [(k1, v1)] → [(k2, v2)]       -- input/output
mkMapper2 f_setup f_map f_cleanup = f_cleanup ∘ foldl f_map f_setup
```

- A state-monadic map (in paper)



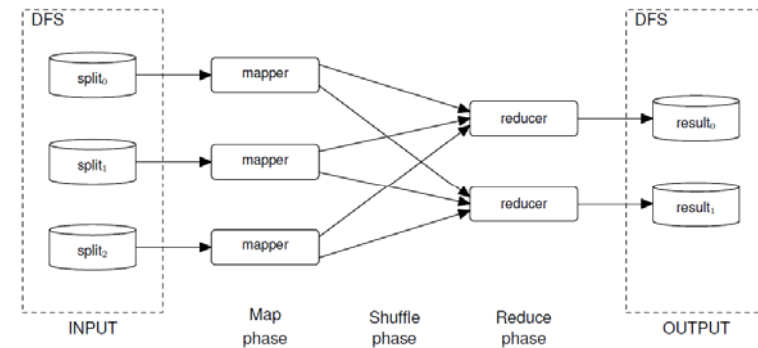
Shuffle Phase in Hadoop





Shuffle Phase in Hadoop

- 3 sub-phases
 1. Partitioning (cf. reduce job)
 2. Sorting
 3. Grouping (cf. reducer func.)



```
shuffleMR :: ((k2, v2) → Integer) -- hashP  
           → (k2 → k2 → Ordering) -- compS  
           → (k2 → k2 → Ordering) -- compG  
           → Bag[(k2, v2)] -- input  
           → Bag[(k2, [v2])] -- output
```

```
shuffleMR hashP compS compG input  
= let aftP = grpByID $ mapB (applyW hashP) $ flattenB input  
    aftS = mapB (sortByKey compS) aftP  
  in mapB (grpByKey compG) aftS
```



Toward High-Level Model

- 3-phase Implementation

1. Partitioning (cf. reduce job)
2. Sorting
3. Grouping (cf. reducer func.)

Sorting and grouping should be consistent



- Possibly any comparator

Let $\text{comp } a \ b = |a - b| < 2$
then $[3, 4, 5, 7] \rightarrow [3, 4, 5], [7]$



- Sorting after grouping

1. Partitioning (cf. reduce job)
2. Grouping (cf. reducer func.)
3. Sorting

cf. secondary-sorting

- Explicit keys

Use (k_P, k_G, k_S) for keys



Difference from Low-Level Model

Low-Level
Model

```
shuffleMR :: ((k2, v2) → Integer)           -- hashP
            → (k2 → k2 → Ordering)         -- compS
            → (k2 → k2 → Ordering)         -- compG
            → Bag[(k2, v2)]                 -- input
            → Bag[(k2, [v2])]               -- output
```



Comparator → Type Class

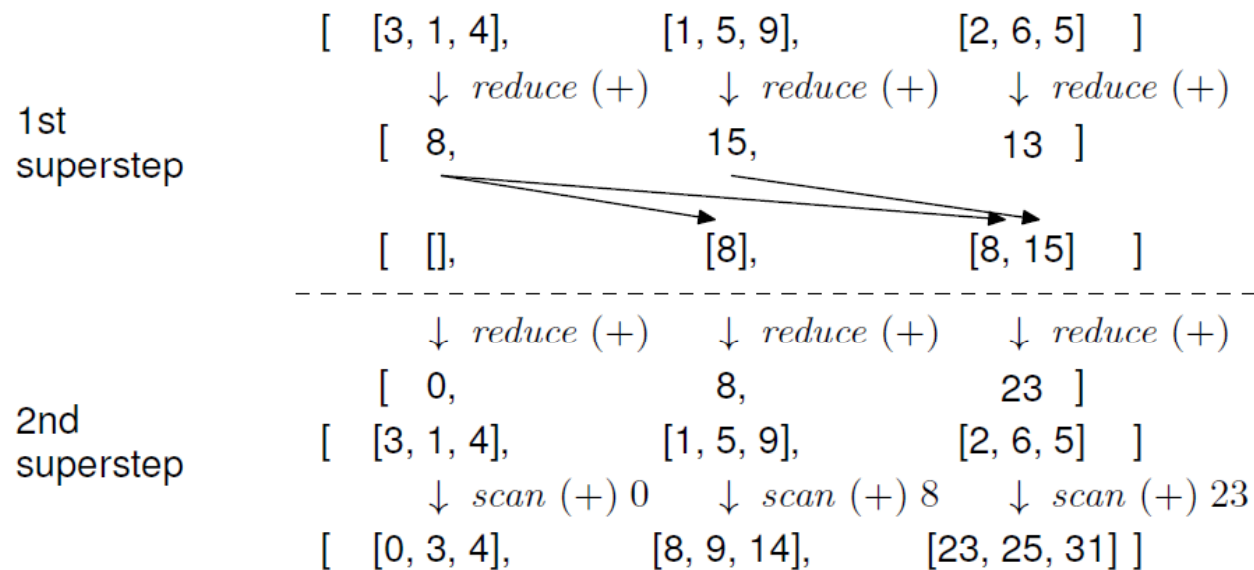
High-Level
Model

```
shuffleMR2 :: Eq kP ⇒ Ord kG ⇒ Ord kS -- (requirements)
            ⇒ Bag [(kP, kG, kS), v2]    -- input
            → Bag [(kP, kG, kS), [v2]] -- output
```

Explicit 3-keys



Case Study: BSP-based Scan



- Challenge

A sequential program → A MapReduce program



A Program for BSP-scan

- Specification: based on **nested-lists**

```
scanBSP :: (a → a → a) → a → Int → [[a]] → [[a]]
scanBSP (⊕) e p xss
  = let ys = map (foldl (⊕) ι⊕) xss
      zss = map (λp. take p ys) [0..(p - 1)]
  in map scan' (zip xss zss)
  where scan' (xs, zs) = scan (⊕) (foldl (⊕) e zs) xs
```

communication

-- 1st superstep

-- 2nd superstep



A MapReduce Program for BSP-scan

- 1-pass MapReduce with state-monadic Mapper

$scanBSPMR (\oplus) e pp xb$

$= mapReduceH mapper reducer xb$

where

$mapper = mkMapper3 f_{setup} f_{map} f_{cleanup}$

$f_{setup} = (0, \iota_{\oplus})$

$f_{map} (-, s) (k, v) = \mathbf{let} p = \mathit{div} k 3$
 $\mathbf{in} ((p, s \oplus v), [((p, 1, k), v)])$

$f_{cleanup} (p, s) = \mathit{map} (\backslash p' \rightarrow ((p', 1, -1000 + p), s)) [p + 1..pp - 1]$

$reducer = mkReducer1 f_{reduce}$

$f_{reduce} ((p, -, -), vs) = \mathbf{let} zs = \mathit{take} p vs$

$xs = \mathit{drop} p vs$

Separate input-list

$\mathbf{in} \mathit{zip} [(3 * p)..(3 * p + 2)]$
 $(scan (\oplus) (\mathit{foldl} (\oplus) e zs) xs)$

Send input-list
to reducer

Design 3-keys (indices)



Related Work

- K. Ono et al.: *SEFM 2011*.
- F. Jiang et al.: *IEICE Japanese Journal 2014*.
 - Proof with Coq system and code extraction
- J. Dörre, et al.: *Concurrency and Computation: Practice and Experience, 2014*.
 - Developed a cost model and performance tuning



Conclusion

- Summary
 - Two functional models of Hadoop MapReduce
 - Low level: based in the implementation of Hadoop
 - High level: User-friendly Shuffle phase
 - Scan (prefix-sum) algorithms on MapReduce
- Future Work
 - Coq library with code generation
 - Comparison of two scan algorithms
 - Functional models for big-graph processing (Pregel)