

An Executable Sequential Specification for Spark Aggregation

Yu-Fang Chen¹, Chih-Duo Hong¹, Ondřej Lengál^{1,2},
Shin-Cheng Mu¹, Nishant Sinha³, Bow-Yaw Wang¹

¹ Academia Sinica, Taiwan

² Brno University of Technology, Czech Republic

³ IBM Research, India

Abstract. Spark is a new promising platform for scalable data-parallel computation. It provides several high-level application programming interfaces (APIs) to perform parallel data aggregation. Since execution of parallel aggregation in Spark is inherently non-deterministic, a natural requirement for Spark programs is to give the same result for any execution on the same data set. We present PURESPARK, an executable formal Haskell specification for Spark aggregate combinators. Our specification allows us to deduce the precise condition for deterministic outcomes from Spark aggregation. We report case studies analyzing deterministic outcomes and correctness of Spark programs.

1 Introduction

Spark [29,1,30] is a popular platform for scalable distributed data-parallel computation based on a flexible programming environment with concise and high-level APIs. Spark is by many considered as the successor of MapReduce [15,25]. Despite its fame, the precursory computational model of MapReduce suffers from I/O congestion and limited programming support for distributed problem solving. Notably, Spark has the following advantages over MapReduce. First, it has high performance due to distributed, cached, and in-memory computation. Second, the platform adopts a relaxed fault tolerant model where sub-results are recomputed upon faults rather than aggressively stored. Third, lazy evaluation semantics is used to avoid unnecessary computation. Finally, Spark offers greater programming flexibility through its powerful APIs founded in functional programming. Spark also owes its popularity to a unified framework for efficient graph, streaming, and SQL-based relational database computation, a machine learning library, and the support of multiple distributed data storage formats. Spark is one of the most active open-source projects with over 1000 contributors [1].

In a typical Spark program, a sequence of transformations followed by an action are performed on Resilient Distributed Datasets (RDDs). An RDD is the principal abstraction for data-parallel computation in Spark. It represents a read-only collection of data items partitioned and stored distributively. RDD operations such as *map*, *reduce*, and *aggregate* are called *combinators*. They generate and aggregate data in RDDs to carry out Spark computation. For instance, the *aggregate* combinator takes user-defined functions *seq* and *comb*: *seq* accumulates a sub-result for each partition while *comb* merges sub-results across different partitions. Spark also provides a family of aggregate combinators for common data structures such as pairs and graphs. In Spark computation, data aggregation is ubiquitous.

Programming in Spark, however, can be tricky. Since sub-results are computed using multiple applications of *seq* and *comb* across partitions concurrently, the order of their applications varies on different executions. Because of indefinite orders of computation, aggregation in Spark is inherently *non-deterministic*. A Spark program may produce different outcomes for the same input on different runs. This form of non-deterministic computation has other side effects. For instance, the private function `AreaUnderCurve.of` in the Spark machine learning library computes numerical integration distributively; it exhibits numerical instability due to non-deterministic computation. Consider the integral of x^{73} on the interval $[-2, 2]$. Since x^{73} is an odd function, the integral is 0. In our experiments, `AreaUnderCurve.of` returns different results ranging from -8192.0 to 12288.0 on the same input because of different orders of floating-point computation. To ensure deterministic outcomes, programmers must carefully develop their programs to adhere to Spark requirements.

Unfortunately, Spark’s documentation does not specify the requirements formally. It only describes informal algebraic properties about combinators to ensure correctness. The documentation provides little help to a programmer in understanding the complex, and sometimes unexpected, interaction between *seq* and *comb*, especially when these two are functions over more complex domains, e.g. lists or trees. Inspecting the Spark implementation is a laborious job since public combinators are built by composing a long chain of generic private combinators—determining the execution semantics from the complex implementation is hard. Moreover, Spark is continuously evolving and the implementation semantics may change significantly across releases. We therefore believe that a formal specification of Spark combinators is necessary to help developers understand the program semantics better, clarify hidden assumptions about RDDs, and help to reason about correctness and sources of non-determinism in Spark programs.

Building a formal specification for Spark is far from straightforward. Spark is implemented in Scala and provides high-level APIs also in Python and Java. Because Spark heavily exploits various language features of Scala, it is hard to derive specifications without formalizing the operational semantics of the Scala language, which is not an easy task by itself. Instead of that, we have developed a Haskell library `PURESPARK` [4], which for each key Spark combinator provides an abstract sequential functional specification in Haskell. We use Haskell as a specification language for two reasons. First, the core of Haskell has strong formal foundations in λ -calculus. Second, program evaluation in Haskell, like in Scala, is lazy, which admits faithful modeling of Spark aggregation. Through the use of Haskell we obtain a concise formal functional model for Spark combinators without formalizing Scala.

An important goal of our specification is to make non-determinism in various combinators explicit. Spark developers can inspect it to identify sources of non-determinism when program executions yield unexpected outputs. Researchers can also use it to understand distributed Spark aggregation and investigate its computational pattern. Our specification is also *executable*. A programmer can use the Haskell APIs to implement data-parallel programs, test them on different input RDDs, and verify correctness of outputs independent of the Spark programming environment. In our case studies, we capture non-deterministic behaviors of real Spark programs by executing the corresponding `PURESPARK` specifications with crafted input data sets. We also show that the sequential specification is useful in developing distributed Spark programs.

Our main contributions are summarized below:

- We present formal, functional, sequential specifications for key Spark aggregate combinators. The PURESPARK specification consists of executable library APIs. It can assist Spark program development by mimicking data-parallel programming in conventional environments.
- Based on the specification, we investigate and identify necessary and sufficient conditions for Spark aggregate combinators to produce deterministic outcomes for general and pair RDDs.
- Our specification allows to deduce the precise condition for deterministic outcomes from Spark aggregation.
- We perform a series of case studies on practical Spark programs to validate our formalization. With PURESPARK, we find instances of numerical instability in the Spark machine learning library.
- Up to our knowledge, this is the first work to provide a formal, functional specification of key Spark aggregate combinators for data-parallel computation.

2 Preliminaries

Let A be a non-empty set and $\odot : A \times A \rightarrow A$ be a function. An element $i \in A$ is the *identity* of \odot if for every $a \in A$, it holds that $a = i \odot a = a \odot i$. The function \odot is *associative* if for every $a, a', a'' \in A$, $a \odot (a' \odot a'') = (a \odot a') \odot a''$; \odot is *commutative* if for every $a, a' \in A$, $a \odot a' = a' \odot a$. The algebraic structure (A, \odot) is a *semigroup* if \odot is associative. A *monoid* is a structure (A, \odot, \perp) such that (A, \odot) is a semigroup and $\perp \in A$ is the identity of \odot . The semigroup (A, \odot) and monoid (A, \odot, \perp) are commutative if \odot is commutative.

Haskell is a strongly typed purely functional programming language. Similar to Scala, Haskell programs are lazily evaluated. We use several widely used Haskell functions (Figure 1). **fst** and **snd** are projections on pairs. **null** tests whether a list is empty. **elem** is the membership function for lists; its infix notation is often used, as in `0 `elem` []`. **(++)** concatenates two lists; it is used as an infix operator, as in `[False] ++ [True]`. **map** applies a function to elements of a list. **reduce!** merges elements of a list by a given binary function from left to right. **fold!** accumulates by applying a function to elements of a list iteratively, also from left to right. **concat** concatenates elements in a list. **concatMap** applies a function to elements of a list and concatenates the results. **lookup** finds the value of a key in a list of pairs. **filter** selects elements from a list by a predicate.

In order to formalize non-determinism in distributed aggregation, we define the following non-deterministic shuffle function for lists:

```
shuffle! :: [ $\alpha$ ]  $\rightarrow$  [ $\alpha$ ]
shuffle! xs = ...      -- shuffle xs randomly
```

A random monad can be used to define random shuffling. Instead of explicit monadic notation, we introduce the *chaotic shuffle!* function in our presentation for the sake of brevity. Thus, **shuffle!** `[0, 1, 2]` evaluates to one of the six possible lists `[0, 1, 2]`, `[0, 2, 1]`, `[1, 0, 2]`, `[1, 2, 0]`, `[2, 0, 1]`, or `[2, 1, 0]` randomly. Using **shuffle!**, more chaotic functions are defined.

```
map! :: ( $\alpha \rightarrow \beta$ )  $\rightarrow$  [ $\alpha$ ]  $\rightarrow$  [ $\beta$ ]
map! f xs = shuffle! (map f xs)

concatMap! :: ( $\alpha \rightarrow [\beta]$ )  $\rightarrow$  [ $\alpha$ ]  $\rightarrow$  [ $\beta$ ]
concatMap! f xs = concat (map! f xs)
```

```

fst :: (α, β) → α
fst (x, _) = x

null :: [α] → Bool
null [] = True
null (x:xs) = False

(++) :: [α] → [α] → [α]
[] ++ ys = ys
x:xs ++ ys = x:(xs ++ ys)

reducel :: (α → α → α) → [α] → α
reducel h (x:xs) = foldl h x xs

concat :: [[α]] → [α]
concat [] = []
concat (xs:xss) = xs ++ (concat xss)

lookup :: α → [(α, β)] → Maybe β
lookup k [] = Nothing
lookup k ((x, y):xys) = if k == x
    then Just y else lookup k xys

snd :: (α, β) → β
snd (_, y) = y

elem :: α → [α] → Bool
elem x [] = False
elem x (y:ys) = x == y || elem x ys

map :: (α → β) → [α] → [β]
map f [] = []
map f (x:xs) = (f x):(map f xs)

foldl :: (β → α → β) → β → [α] → β
foldl h z [] = z
foldl h z (x:xs) = foldl h (h z x) xs

concatMap :: (α → [β]) → [α] → [β]
concatMap xs = concat (map f xs)

filter :: (α → Bool) → [α] → [α]
filter p [] = []
filter p (x:xs) = if p x
    then x:(filter p xs) else filter p xs

```

Chaotic **map!** shuffles the result of **map** randomly, **concatMap!** concatenates the shuffled result of **map**. For instance, **map! even** [0, 1] evaluates to **[False, True]** or **[True, False]**; **concatMap! fact** [2, 3] evaluates to [1, 2, 1, 3] or [1, 3, 1, 2] where **fact** computes a sorted list of factors (note that the two sub-sequences [1,2] and [1,3] are kept intact).

The function **repartition!** shuffles a given list and partitions the shuffled list into several non-empty lists. For instance, **repartition!** [0, 1] results in [[0], [1]], [[1], [0]], [[0, 1]], or [[1, 0]]. The chaotic function can be implemented by a random monad easily; its precise definition is omitted here.

Resilient Distributed Datasets (RDDs) are the basic data abstraction in Spark. An RDD is a collection of partitions of immutable data; data in different partitions can be processed concurrently. We formalize partitions by lists, and RDDs by lists of partitions.

The Spark aggregate combinator computes *sub-results* of every partitions in an RDD, and returns the aggregated result by combining sub-results.

More concretely, let z be a default aggregated value. **aggregate** applies **foldl** seq z to every partition of rdd . Hence the sub-result of each partition is accumulated by folding elements in the partition with **seq**. The combinator then combines sub-results by another folding using **comb**.

Note that the chaotic **map!** function is used to model non-deterministic interleavings of sub-results. To exploit concurrency, Spark creates a task to compute the sub-result for each partition. These tasks are executed concurrently and hence induce non-deterministic computation. We use the chaotic **map!** function to designate non-determinism explicitly.

A related combinator is **reduce**. Instead of **foldl**, the combinator uses **reducel** to aggregate data in an RDD.

```
reduce :: ( $\alpha \rightarrow \alpha \rightarrow \alpha$ )  $\rightarrow$  RDD  $\alpha \rightarrow \alpha$ 
reduce comb rdd = let presults = map! (reducel comb) rdd
                  in reducel comb presults
```

Similar to the **aggregate** combinator, **reduce** computes sub-results concurrently. The chaotic **map!** function is again used to model non-deterministic computation.

Sub-results of different partitions are computed in parallel, but the **aggregate** combinator still combines sub-results sequentially. This can be further parallelized. Observe that several sub-results may be available simultaneously from distributed computation. The Spark **treeAggregate** combinator applies **comb** to pairs of sub-results concurrently until the final result is obtained. In addition to concurrent computation of sub-results, **treeAggregate** also combines sub-results from different partitions in parallel.

In our specification, two chaotic functions are used to model non-deterministic computation on two different levels. The **map!** function models non-determinism in computing sub-results of partitions. The **apply!** function (introduced below) models concurrent combination of sub-results from different partitions. It combines two consecutive sub-results picked chaotically, and repeats such chaotic combinations until the final result is obtained. Observe that the computation has a binary-tree structure with **comb** as internal nodes and sub-results from different partitions as leaves.

```
apply! :: ( $\beta \rightarrow \beta \rightarrow \beta$ )  $\rightarrow$  [ $\beta$ ]  $\rightarrow \beta$ 
apply! comb [r] = r
apply! comb [r, r'] = comb r r'
apply! comb rs = let (ls', l', r', rs') = ...    --  $rs == ls' \uparrow l' \uparrow r' \uparrow rs'$ 
                  in apply! comb (ls'  $\uparrow$  [comb l' r']  $\uparrow$  rs')
```

```
treeAggregate ::  $\beta \rightarrow (\beta \rightarrow \alpha \rightarrow \beta) \rightarrow (\beta \rightarrow \beta \rightarrow \beta) \rightarrow$  RDD  $\alpha \rightarrow \beta$ 
treeAggregate z seq comb rdd = let presults = map! (foldl seq z) rdd
                              in apply! comb presults
```

The **treeReduce** combinator optimizes **reduce** by combining sub-results in parallel. Similar to **treeAggregate**, two levels of non-deterministic computation can occur.

```
treeReduce :: ( $\alpha \rightarrow \alpha \rightarrow \alpha$ )  $\rightarrow$  RDD  $\alpha \rightarrow \alpha$ 
treeReduce comb rdd = let presults = map! (reducel comb) rdd
                      in apply! comb presults
```

Pair RDDs. Key-value pairs are widely used in data parallel computation. If the data type of an RDD is a pair, we say that the RDD is a *pair* RDD. The first and second elements in a pair are called the *key* and the *value* of the pair respectively.

type PairRDD $\alpha \beta = \text{RDD } (\alpha, \beta)$

In a pair RDD, different pairs can have the same key. Spark provides combinators to aggregate values associated with the same key. The **aggregateByKey** combinator returns an RDD by aggregating values associated with the same key. We use the following functions to formalize **aggregateByKey**:

hasKey :: $\alpha \rightarrow \text{Partition } (\alpha, \beta) \rightarrow \text{Bool}$	hasValue :: $\alpha \rightarrow \beta \rightarrow \text{Partition } (\alpha, \beta) \rightarrow \beta$
hasKey k ps = case (lookup k ps) of	hasValue k val ps = case (lookup k ps) of
Just _ \rightarrow True	Just v \rightarrow v
Nothing \rightarrow False	Nothing \rightarrow val

addTo :: $\alpha \rightarrow \beta \rightarrow \text{Partition } (\alpha, \beta) \rightarrow \text{Partition } (\alpha, \beta)$
addTo key val ps = **foldl** ($\lambda r (k, v) \rightarrow$ **if** key == k **then** r **else** (k, v):r) [(key, val)] ps

The expression **hasKey** k ps checks if key appears in a partition of pairs. **hasValue** k val ps finds a value associated with key in a partition of pairs. It evaluates to the default value val if key does not appear in the partition. The expression **addTo** key val ps adds the pair (key, val) to the partition ps, and removes other pairs with the same key.

The **aggregateByKey** combinator first aggregates all pairs with the value z and the function mergeComb in each partition. If values vs are associated with the same key in a partition, the value **foldl** mergeComb z vs for the key is pre-aggregated. Since a key may appear in several partitions, all pre-aggregated values associated with the key across different partitions are merged using mergeValue.

aggregateByKey :: $\gamma \rightarrow (\gamma \rightarrow \beta \rightarrow \gamma) \rightarrow (\gamma \rightarrow \gamma \rightarrow \gamma) \rightarrow \text{PairRDD } \alpha \beta \rightarrow \text{PairRDD } \alpha \gamma$
aggregateByKey z mergeComb mergeValue pairRdd =
 let mergeBy fun left (k, v) = **addTo** k (fun (hasValue k z left) v) left
 preAgg = **concatMap!** (**foldl** (mergeBy mergeComb) []) pairRdd
 in **repartition!** (**foldl** (mergeBy mergeValue) [] preAgg)

In the specification, we accumulate values associated with the same key by mergeComb in each partition, keeping a list of pairs of a key and the partially aggregated value for the key. Since accumulation in different partitions runs in parallel, the chaotic **concatMap!** function is used to model such non-deterministic computation. After all partitions finish their accumulation, mergeValue merges values associated with the same key across different partitions. The final pair RDD can have a default or user-defined partitioning. Since a user-defined partitioning may shuffle a pair RDD arbitrarily, it is in our specification modeled by the chaotic **repartition!** function.

Pair RDDs have a combinator corresponding to **reduce** called **reduceByKey**. **reduceByKey** merges all values associated with a key by mergeValue, following a similar computational pattern as **aggregateByKey**. Note that every key is associated with at most one value in resultant pair RDDs of **aggregateByKey** or **reduceByKey**.

reduceByKey :: $(\beta \rightarrow \beta \rightarrow \beta) \rightarrow \text{PairRDD } \alpha \beta \rightarrow \text{PairRDD } \alpha \beta$
reduceByKey mergeValue pairRdd =
 let merge left (k, v) = **case** lookup k left of **Just** v' \rightarrow **addTo** k (mergeValue v' v) left
 Nothing \rightarrow **addTo** k v left
 preAgg = **concatMap!** (**foldl** merge []) pairRdd
 in **repartition!** (**foldl** merge [] preAgg)

Spark also provides a library, called GraphX, for a distributed analysis of graphs. See [12] for a formalization of some of its key functions.

4 Deterministic Aggregation

Having deterministic outcomes is desired from all aggregation functions. If a function may return different values on different executions, the function is often not implemented correctly. A program with explicit assumptions on the input data is also desirable. Otherwise, the program may work correctly on certain data sets but produce unexpected outcomes on others where implicit assumptions do not hold [27]. We now investigate conditions under which Spark aggregation combinators always produce deterministic outcomes. Proofs of the given lemmas can be found in [12]. Proofs of some crucial lemmas have also been formalized using Agda [4].

We first show how to deal with non-deterministic behaviors in the **aggregate** combinator. Consider a variant of the formalization of **aggregate** from Section 3:

```
aggregate' ::  $\beta \rightarrow (\beta \rightarrow \alpha \rightarrow \beta) \rightarrow (\beta \rightarrow \beta \rightarrow \beta) \rightarrow \text{RDD } \alpha \rightarrow \beta$ 
aggregate' z seq comb rdd = let presalts = perm (map (foldl seq z) rdd)
  in foldl comb z presalts
```

Observe that we changed the application of the chaotic **map!** function with an application of the permutation perm after the regular **map** function. The function composition perm(**map** ...) is a concrete instantiation of **map!**, that is, a function that permutes its list argument. Notice that perm can be pushed inside **map**:

```
perm (map f xs) == map f (perm xs).
```

Assume that rdd was obtained from a list xs by splitting and permuting, that is, $\text{rdd} == \text{perm}' (\text{split } xs)$ where $\text{split} :: [\alpha] \rightarrow [[\alpha]]$ satisfies $xs == (\text{concat} \cdot \text{split}) xs$. We can therefore rewrite the computation of presalts in **aggregate**['] to

```
let pres = perm (map (foldl seq z) (perm' (split xs))),
```

After pushing perm inside map, we obtain

```
let pres = map (foldl seq z) ((perm . perm') (split xs)).
```

Since perm . perm' is also a permutation perm'', we have

```
let pres = map (foldl seq z) rdd'
```

where rdd' is another RDD obtained from xs by splitting and shuffling. Let us call (deterministic) instances of **repartition!** as *partitionings*. As a consequence, we focus only on proving if calls to **aggregate**^D defined below have deterministic outcomes for different partitionings of a list into RDDs:

```
aggregateD ::  $\beta \rightarrow (\beta \rightarrow \alpha \rightarrow \beta) \rightarrow (\beta \rightarrow \beta \rightarrow \beta) \rightarrow \text{RDD } \alpha \rightarrow \beta$ 
aggregateD z seq comb rdd = let pres = map (foldl seq z) rdd
  in foldl comb z pres
```

Moreover, we define deterministic versions of **reduce**

```
reduceD ::  $(\alpha \rightarrow \alpha \rightarrow \alpha) \rightarrow \text{RDD } \alpha \rightarrow \alpha$ 
reduceD comb rdd = let presalts = perm (map (reduce! comb) rdd)
  in reduce! comb presalts
```

and also **treeAggregate**^D and **treeReduce**^D in a similar way.

In the following, given a function f that takes an RDD as one of its parameters and contains a single occurrence of the chaotic **map!** (respectively **concatMap!**) function, we use f^D to denote the function obtained from f by replacing the chaotic **map!** (respectively **concatMap!**) with a regular **map** (respectively **concatMap**). A similar reasoning

can show that it suffices to check whether calls to f^D have deterministic outcomes for different partitionings on a list into RDDs.

For better readability, standard mathematical notation of functions is used in the rest of this section. We represent a Haskell function application $f\ x_1 \dots x_n$ as $f(x_1, \dots, x_n)$.

4.1 aggregate

In this section, we give conditions for deterministic outcomes of calls to the aggregate combinator **aggregate**(z, seq, \oplus, rdd) for $z :: \beta$, $seq :: \beta \times \alpha \rightarrow \beta$, $\oplus :: \beta \times \beta \rightarrow \beta$, and $rdd :: \text{RDD } \alpha$. We first define what it means for calls to the **aggregate** combinator to have deterministic outcomes.

Definition 1. *Calls to **aggregate**(z, seq, \oplus, rdd) have deterministic outcomes if*

$$\mathbf{aggregate}^D(z, seq, \oplus, \text{part}(L)) = \mathbf{foldl}(seq, z, L) \quad (1)$$

for all lists L and partitionings part .

Conventionally, **aggregate** is regarded as a parallelized counterpart of **foldl**. For example, the sequential **aggregate** function in the standard Scala library ignores the \oplus operator and is implemented by **foldl**. This is why we characterize deterministic **aggregate** as **foldl** in Definition 1. Our characterization, however, does not cover all **aggregate** calls that always give the same outputs. In particular, it does not cover an **aggregate** call where \oplus is a constant function, which is, however, quite suspicious in a distributed data-parallel computation and should be reported.

We give necessary and sufficient conditions for **aggregate** calls to have deterministic outcomes in several lemmas, culminating in Corollary 1. The first lemma allows us to check only conditions on seq and \oplus over all possible pairs of lists instead of enumerating all possible partitionings on lists. For brevity, we use $\langle p_1 \rangle$ for $\mathbf{foldl}(seq, z, p_1)$, and $\text{img}(\mathbf{foldl}(seq, z))$ for the image of $\mathbf{foldl}(seq, z, L)$ for any list L . That is, $\text{img}(\mathbf{foldl}(seq, z)) = \{y \mid \text{there is a list } L \text{ such that } \mathbf{foldl}(seq, z, L) = y\}$.

Lemma 1. *Calls to **aggregate**(z, seq, \oplus, rdd) have deterministic outcomes iff:*

1. *($\text{img}(\mathbf{foldl}(seq, z)), \oplus, z$) is a commutative monoid, and*
2. *for all lists $p_1, p_2 :: [\alpha]$, $\langle p_1 \uplus p_2 \rangle = \langle p_1 \rangle \oplus \langle p_2 \rangle$.*

Note that condition 2 in Lemma 1 is equivalent to saying that $\langle \cdot \rangle$ is a list homomorphism to the monoid $(\text{img}(\mathbf{foldl}(seq, z)), \oplus, z)$ [6].

The lemma below further helps us reduce the need of testing conditions over all possible pairs of lists to conditions over elements of $\alpha \times \text{img}(\mathbf{foldl}(seq, z))$.

Lemma 2. *Let \oplus be associative on $\gamma = \text{img}(\mathbf{foldl}(seq, z))$ and z be the identity of \oplus on γ . The following are equivalent:*

1. *for all lists $p_1, p_2 :: [\alpha]$,*

$$\langle p_1 \uplus p_2 \rangle = \langle p_1 \rangle \oplus \langle p_2 \rangle, \quad (2)$$
2. *for all elements $d :: \alpha$ and $e :: \gamma$,*

$$seq(e, d) = e \oplus seq(z, d). \quad (3)$$

Summarizing the lemmas, we get the following corollary:

Corollary 1. *Calls to **aggregate**(z, seq, \oplus, rdd) have deterministic outcomes iff*

1. *($\text{img}(\mathbf{foldl}(seq, z)), \oplus, z$) is a commutative monoid and*
2. *for all $d :: \alpha$ and $e :: \text{img}(\mathbf{foldl}(seq, z))$, it holds that $seq(e, d) = e \oplus seq(z, d)$.*

4.2 reduce

This section explores conditions for deterministic outcomes of calls to **reduce**(\oplus, rdd) for $\oplus :: \alpha \times \alpha \rightarrow \alpha$ and $rdd :: \text{RDD } \alpha$. We use the function **reduce**^D defined in the introduction of Section 4. For **reduce**, we assume that for any non-empty list, all partitions of its partitioning are non-empty (otherwise the result of **reduce** is undefined).

We define deterministic outcomes for **reduce** as follows.

Definition 2. Calls to **reduce**(\oplus, rdd) have deterministic outcomes if

$$\mathbf{reduce}^D(\oplus, \text{part}(L)) = \mathbf{reducel}(\oplus, L) \quad (4)$$

for all lists L and partitionings part .

We reduce the problem of checking if **reduce** has deterministic outcomes to the problem of checking if **aggregate** has deterministic outcomes by the following lemma.

Lemma 3. Calls to **reduce**(\oplus, rdd) have deterministic outcomes iff calls to **aggregate**(**Nothing**, seq' , \oplus' , rdd) have deterministic outcomes, where seq' and \oplus' are as follows:

$$\begin{array}{ll} \text{seq}' \times y = \text{case } x \text{ of} & (\oplus') \times y = \text{case } (x, y) \text{ of } (\text{Nothing}, y') \rightarrow y' \\ \text{Nothing} \rightarrow \text{Just } y & (x', \text{Nothing}) \rightarrow x' \\ \text{Just } x' \rightarrow \text{Just } (x' \oplus y) & (\text{Just } x', \text{Just } y') \rightarrow \text{Just } (x' \oplus y'). \end{array}$$

Combining Corollary 1 and Lemma 3, we get the condition for deterministic outcomes of **reduce**(\oplus, rdd) calls.

Corollary 2. Calls to **reduce**(\oplus, rdd) have deterministic outcomes iff (α, \oplus) is a commutative semigroup.

4.3 treeAggregate and treeReduce

This section gives conditions for deterministic outcomes of calls to the following two aggregate combinators:

1. **treeAggregate**($z, \text{seq}, \oplus, rdd$) for $z :: \beta$, $\text{seq} :: \beta \times \alpha \rightarrow \beta$, $\oplus :: \beta \times \beta \rightarrow \beta$, and $rdd :: \text{RDD } \alpha$; and
2. **treeReduce**(\oplus, rdd) for $\oplus :: \alpha \times \alpha \rightarrow \alpha$, $rdd :: \text{RDD } \alpha$.

Different from **aggregate** and **reduce**, the tree variants have another level of non-determinism modeled by **apply!**. The chaotic function effectively simulates non-deterministic computation with a binary-tree structure (Section 3).

To define calls to **treeAggregate** and **treeReduce** to have deterministic outcomes, we use the functions **treeAggregate**^T and **treeReduce**^T obtained by adding an explicit deterministic instantiation of **apply!** to **treeAggregate**^D and **treeReduce**^D.

Definition 3. Calls to **treeAggregate**($z, \text{seq}, \oplus, rdd$) and **treeReduce**(\oplus, rdd) have deterministic outcomes if

$$\mathbf{treeAggregate}^T(\text{apply}, z, \text{seq}, \oplus, \text{part}(L)) = \mathbf{foldl}(\text{seq}, z, L) \quad (5)$$

and

$$\mathbf{treeReduce}^T(\text{apply}, \oplus, \text{part}(L)) = \mathbf{reducel}(\oplus, L) \quad (6)$$

respectively for all lists L , partitionings part , and instantiations apply of **apply!**.

The following two propositions state necessary and sufficient conditions for the **treeAggregate** and **treeReduce** combinators to have deterministic outcomes.

Proposition 1. Calls to **treeAggregate**($z, \text{seq}, \oplus, rdd$) have deterministic outcomes iff calls to **aggregate**($z, \text{seq}, \oplus, rdd$) have deterministic outcomes.

Proposition 2. Calls to **treeReduce**(\oplus, rdd) have deterministic outcomes iff calls to **reduce**(\oplus, rdd) have deterministic outcomes.

4.4 aggregateByKey and reduceByKey

We proceed by investigating conditions for the following combinators on pair RDDs:

1. **aggregateByKey**($z, seq, \oplus, prdd$) for $z :: \gamma$, $seq :: \gamma \times \beta \rightarrow \gamma$, $\oplus :: \gamma \times \gamma \rightarrow \gamma$, and $prdd :: \text{PairRDD } \alpha \beta$; and
2. **reduceByKey**($\oplus, prdd$) for $\oplus :: \beta \times \beta \rightarrow \beta$ and $prdd :: \text{PairRDD } \alpha \beta$.

We define an auxiliary function **filterkey** that obtains a list of all values associated with the given key from a list of pairs.

filterkey $:: \alpha \rightarrow [(\alpha, \beta)] \rightarrow [\beta]$

filterkey $_ [] = []$

filterkey $k (k, v):xs = v:(\text{filterkey } k \text{ } xs)$

filterkey $k (_, _):xs = \text{filterkey } k \text{ } xs$

Deterministic outcomes of calls to **aggregateByKey** are now defined using the function **aggregateByKey**^D as follows.

Definition 4. Calls to **aggregateByKey**($z, seq, \oplus, prdd$) have deterministic outcomes if

$$\text{lookup}(k, \text{aggregateByKey}^D(z, seq, \oplus, \text{part}(L))) = \text{foldl}(z, seq, \text{filterkey}(k, L))$$

for all lists L of pairs, partitionings part , and keys k .

Finally, the following proposition states the conditions that need to hold for calls to **aggregateByKey** to have deterministic outcomes.

Proposition 3. Calls to **aggregateByKey**($z, seq, \oplus, prdd$) have deterministic outcomes iff calls to **aggregate**(z, seq, \oplus, rdd) have deterministic outcomes.

We define when calls to **reduceByKey** have deterministic outcomes via **reduceByKey**^D.

Definition 5. Calls to **reduceByKey**($\oplus, prdd$) have deterministic outcomes if

$$\text{lookup}(k, \text{reduceByKey}^D(\oplus, \text{part}(L))) = \text{reducel}(\oplus, \text{filterkey}(k, L))$$

for all list L of pairs, partitioning part , and key k .

Proposition 4. Calls to **reduceByKey**($\oplus, prdd$) have deterministic outcomes iff calls to **reduce**(\oplus, rdd) have deterministic outcomes.

4.5 Discussion

Our conditions for deterministic outcomes are more general than it appears. In addition to scalar data, such as integers, they are also applicable to RDDs containing non-scalar data, such as lists or sets. In our extended set of case studies, we will prove deterministic outcomes from a distributed Spark program using non-scalar data [12].

Corollary 1 gives necessary and sufficient conditions for calls to **aggregate** to have deterministic outcomes. Instead of checking whether **aggregate** computes the same

result on all possible partitionings on any list for given z , seq , and $comb$, the corollary, instead, allows us to investigate properties for all elements of $img(\mathbf{foldl}(seq, z)) \times img(\mathbf{foldl}(seq, z))$ and $\alpha \times img(\mathbf{foldl}(seq, z))$. Our precise conditions reduce the need of checking all partitionings to checking all elements of Cartesian products. It appears that deterministic outcomes from calls to combinators can be verified automatically. The problem, however, remains difficult for the following reasons:

- (a) The domain $img(\mathbf{foldl}(seq, z))$ can be infinite and in general not computable.
- (b) Even if α and $img(\mathbf{foldl}(seq, z))$ are computable, seq and \oplus may not be computable. Naïvely enumerating elements in α and $img(\mathbf{foldl}(seq, z))$ would not work.
- (c) Testing equality between elements of $img(\mathbf{foldl}(seq, z))$ can be undecidable.

Given $seq :: \beta \times \alpha \rightarrow \beta$, recall that $img(\mathbf{foldl}(seq, z))$ is a subset of β . A sound but incomplete way to avoid (a) in practice is to test the properties of \oplus on all elements of β instead. If a counterexample is found for some elements of β , the counterexample may not be valid in a real **aggregate** call because it may not belong to $img(\mathbf{foldl}(seq, z))$. In practical cases, the sets α and β are finite (such as machine integers) and equality between their elements is decidable. Even for such cases, checking if outcomes of **aggregate** are deterministic is still difficult since seq and \oplus might not terminate for some input. In many real Spark programs, however, seq and \oplus are very simple and thus computable (for instance, with only bounded loops or recursion). A semi-procedure to test these conditions might work on such practical examples.

5 Case Studies

We evaluated advantages of our PURESPARK specification on several case studies. In this section, we first analyze a Spark implementation of linear classification. Using the **treeAggregate** specification and its criteria for deterministic outcomes, we construct inputs yielding non-deterministic outcomes from the Spark implementation. Second, we analyze an implementation of a standard scaler and find a non-deterministic behavior there, too. Yet more case studies are provided in [12].

5.1 Linear Classification

Linear classification is a well-known machine learning technique to classify data sets. Fix a set of *features*. A *data point* is a vector of numerical feature values. A *labeled* data point is a data point with a discrete label. Given a labeled data set, the *classification problem* is to classify (new) unlabeled data points by the labeled data set. A particularly useful subproblem is the *binary* classification problem. Consider, for instance, a data set of vital signs of some population; each data point is labeled by the diagnosis of a disease (positive or negative). The binary classification problem can be used to predict whether a person has the particular disease. Linear classification solves the binary classification problem by finding an optimal hyperplane to divide the labeled data points. After a hyperplane is obtained, linear classification predicts an unlabeled data point by the half-space containing the point. Logistic regression and linear Support Vector Machines (SVMs) are linear classification algorithms.

Consider a data set $\{(\vec{x}_i, y_i) : 1 \leq i \leq n\}$ of data points $\vec{x}_i \in \mathbb{R}^d$ labeled by $y_i \in \{0, 1\}$. Linear classification can be expressed as a numerical optimization problem:

$$\min_{\vec{w} \in \mathbb{R}^d} f(\vec{w}) \quad \text{with} \quad f(\vec{w}) = \xi R(\vec{w}) + \frac{1}{n} \sum_{i=1}^n L(\vec{w}; \vec{x}_i, y_i)$$

where $\xi \geq 0$ is a *regularization parameter*, $R(\vec{w})$ is a *regularizer*, and $L(\vec{w}; \vec{x}_i, y_i)$ is a *loss function*. A vector \vec{w} corresponds to a hyperplane in the data point space. The vector \vec{w}_{opt} attaining the optimum hence classifies unlabeled data points with criteria defined by the objective function $f(\vec{w})$. Logistic regression and linear SVM are but two instances of the optimization problem with objective functions defined by different regularizers and loss functions.

In the Spark machine learning library, the numerical optimization problem is solved by gradient descent. Very roughly, gradient descent finds a local minimum of $f(\vec{w})$ by “walking” in the opposite direction of the gradient of $f(\vec{w})$. The mean of subgradients at data points is needed to compute the gradient of $f(\vec{w})$. The Spark machine learning library invokes **treeAggregate** to compute the mean. Floating-point addition is used as the comb parameter of the aggregate combinator. Since floating-point addition is not associative, we expect to observe non-deterministic outcomes (Proposition 1).

Consider the following three labeled data points: -10^{20} labeled with 1, 600 labeled with 0, and 10^{20} labeled with 1. We create a 20-partition RDD with an equal number of the three labeled data points. The Spark machine learning library function `LogisticRegressionWithSGD.train` is used to generate a logistic regression model to predict the data points -10^{20} , 600, and 10^{20} in each run. Among 49 runs, 19 of them classify the three data points into two different classes: the two positive data points are always classified in the same class, while the negative data point in the other. The other 30 runs, however, classify all three data points into the same class. We observe similar predictions from `SVMWithSGD.train` with the same labeled data points. 37 out of 46 runs classify the data points into two different classes; the other 9 runs classify them into one class. Interestingly, the data points are always classified into two different classes by both logistic regression and linear SVM when the input RDD has only three partitions. As we expected from our analysis of the function, non-deterministic outcomes were witnessed in our Spark distributed environment.

5.2 Standard Scaler

Standardization of data sets is a common pre-processing step in machine learning. Many machine learning algorithms tend to perform better when the training set is similar to the standard normal distribution. In the Spark machine learning library, the class `StandardScaler` is provided to standardize data sets. The function `StandardScaler.fit` takes an RDD of raw data and returns an instance of `StandardScalerModel` to transform data points. Two transformations are available in `StandardScalerModel`. One standardizes a data point by mean, and the other normalizes by variance of raw data. If data points in raw data are transformed by mean, the transformed data points have the mean equal to 0. Similarly, if they are transformed by variance, the transformed data points have the variance 1.

The `StandardScaler` implementation uses **treeAggregate** to compute statistical information. It uses floating-point addition to combine means of raw data in different partitions. As in the previous use case, since floating-point addition is not associative, `StandardScaler` does not produce deterministic outcomes (Section 4.3). In our experiment, we create a 100-partition RDD with values -10^{20} , 600, 10^{20} of the same number

of occurrences. The mean of the data set is $(-10^{20} \times n + 600 \times n + 10^{20} \times n) / (3n) = 200$ where n is the number of occurrences of each value. The data point 200 should therefore be after standardization transformed to 0. In 50 runs on the same data set in our distributed Spark platform, `StandardScaler` transforms 200 to a range of values from -944 to 1142 , validating our prediction of a non-deterministic outcome.

6 Related Work

MapReduce modeling and optimization. In the MapReduce (MR) computation, various cost and performance models have been proposed [25,18,16,31]. These models estimate the execution time and resource requirements of MR jobs. Karloff et al. developed a formal computation model for MR [21] and showed how a variety of algorithms can exploit the combination of sequential and parallel computation in MR. We are not aware of a similar work in the context of Spark. To the best of our knowledge, our work is the first to address the problem of formal, functional specification of Spark aggregation. Verifying the correctness of a MR program involves checking the commutativity and associativity of the reduce function. Xu et al. propose various semantic criteria to model commonly held assumptions on MR programs [28], including determinism, partition isolation, commutativity, and associativity of map/reduce combinators. Their empirical survey shows that these criteria are often overlooked by programmers and violated in practice. A recent survey [27] has found that a large number of industrial MR programs are, in fact, non-commutative. Recent work has proposed techniques for checking commutativity of bounded reducers automatically [13]. Because it is non-trivial to implement high-level algorithms using the MR framework, various approaches to compute optimized MR implementations have been proposed [17,23,24]. Emoto et al. [17] formalize the algebraic conditions using semiring homomorphism, under which an efficient program based on the generate-test-aggregate programming model can be specified in the MR framework. Given a monolithic *reduce* function, the work in [23] tries to decompose *reduce* into partial aggregation functions (similar to *seq* and *comb* in this paper) using program inversion techniques. MOLD [24] translates imperative Java code into MR code by transforming imperative loops into *fold* combinators using semantic-preserving program rewrite rules.

Numerical Stability under MapReduce. Several works try to scale up machine learning algorithms for large datasets using MapReduce [14,25]. To achieve numerically stable results across multiple runs [5,26], for example, preventing overflow, underflow and round-off errors due to finite-precision arithmetic, a variety of techniques are proposed [26]: generalizing sequential numerical stability techniques to distributed settings, shifting data values by constants, divide-and-conquer, etc. We showed that simulating machine learning algorithms using our specification enables early detection of points of numerical instability.

Relational Query Optimization. Relational query optimization is an extensively researched topic [11,20]: the goal is to obtain equivalent but more efficient query expressions by exploiting the algebraic properties of the constituent operators, for instance, join, select, together with statistics on relations and indices. For example, while inner joins commute independent of data, left joins commute only in specific cases. Query optimization for partitioned tables has received less attention [19,2]: because the key relational operators are not partition-aware, most work has focused on necessary but not sufficient conditions for query equivalence. In contrast, we investigate determinism

of Spark aggregate expressions, constructed using partition-aware *seq* and *comb* combinators. We describe necessary and sufficient conditions under which these computations yield deterministic results independent of the data partitions.

Deterministic Parallel Programming. In order to enable deterministic-by-default parallel programming [7,10,8,9,22], researchers have developed several programming abstractions and logical specification languages to ensure that programs produce the same output for the same input independent of thread scheduling. For example, Deterministic Parallel Java [7,8] ensures exclusive writes to shared memory regions by means of verified, user-provided annotations over memory regions. In contrast, deterministic outcomes from Spark aggregation depend on algebraic properties like commutativity and associativity of *seq* and *comb* functions and their interplay

7 Conclusion

In this paper, we give a Haskell specification for various Spark aggregate combinators. We focus on aggregation of RDDs representing general sets, sets of pairs, and graphs. Based on our specification, we derive necessary and sufficient conditions that guarantee deterministic outcomes of the considered Spark aggregate combinators. We investigate several case studies and use the conditions to predict non-deterministic outcomes. Our executable specification can be used by developers for more detailed analysis and efficient development of distributed Spark programs. We also believe that our specifications are valuable resources for research communities to understand Spark better.

There are several future directions. The conditions for deterministic outcomes of aggregate combinators could be used for: (i) creating fully mechanized proofs for properties about data-parallel programs; (ii) developing automatic techniques for detecting non-deterministic outcomes of data-parallel programs; and (iii) synthesizing deterministic concurrent programs from sequential specifications. We have formalized the proofs of some crucial lemmas in Agda [4]. Using Scalaz [3], verified Haskell specifications can be translated to Spark programs to ensure determinism by construction.

Acknowledgement. This work was supported by the Czech Science Foundation (project 17-12465S), the BUT FIT project FIT-S-17-4014, the IT4IXS: IT4Innovations Excellence in Science project (LQ1602), and Ministry of Science and Technology, R.O.C. (MOST projects 103-2221-E-001-019-MY3 and 103-2221-E-001-020-MY3).

References

1. Apache Spark, <https://github.com/apache/spark>
2. IBM DB2 Version 9.7. Partitioned Tables, <https://ibm.biz/BdHyYR>
3. The Scalaz project, <https://github.com/scalaz>
4. PURESPARK, <https://github.com/guluchen/purespark>
5. Bennett, J., Grout, R., Pebay, P., Roe, D., Thompson, D.: Numerically stable, single-pass, parallel statistics algorithms. In: CLUSTER. pp. 1–8 (2009)
6. Bird, R.S.: An introduction to the theory of lists. In: the NATO Advanced Study Institute on Logic of programming and calculi of discrete design. pp. 5–42. Springer (1987)
7. Bocchino, Jr., R.L., Adve, V.S., Dig, D., Adve, S.V., Heumann, S., Komuravelli, R., Overbey, J., Simmons, P., Sung, H., Vakilian, M.: A type and effect system for deterministic parallel Java. In: OOPSLA. pp. 97–116 (2009)
8. Bocchino, Jr., R.L., Heumann, S., Honarmand, N., Adve, S.V., Adve, V.S., Welc, A., Shpeisman, T.: Safe nondeterminism in a deterministic-by-default parallel language. SIGPLAN Not. 46(1), 535–548 (2011)

9. Budimlic, Z., Burke, M.G., Cavé, V., Knobe, K., Lowney, G., Newton, R., Palsberg, J., Peixotto, D.M., Sarkar, V., Schlimbach, F., Tasirlar, S.: Concurrent collections. *Scientific Programming* 18(3-4), 203–217 (2010)
10. Burnim, J., Sen, K.: Asserting and checking determinism for multithreaded programs. *Commun. ACM* 53(6), 97–105 (2010)
11. Chaudhuri, S.: An overview of query optimization in relational systems. *PODS '98* (1998)
12. Chen, Y., Hong, C., Lengál, O., Mu, S., Sinha, N., Wang, B.: An executable sequential specification for Spark aggregation. *arXiv:1702.02439 [cs.DC]* (2017)
13. Chen, Y., Hong, C., Sinha, N., Wang, B.: Commutativity of reducers. In: *Proc. of TACAS'15*. pp. 131–146. LNCS, Springer (2015)
14. Chu, C., Kim, S.K., Lin, Y., Yu, Y., Bradski, G.R., Ng, A.Y., Olukotun, K.: Map-Reduce for machine learning on multicore. In: *NIPS*. pp. 281–288 (2006)
15. Dean, J., Ghemawat, S.: MapReduce: A flexible data processing tool. *Commun. ACM* 53(1), 72–77 (2010)
16. Dörre, J., Apel, S., Lengauer, C.: Modeling and optimizing MapReduce programs. *Concurrency and Computation: Practice and Experience* 27(7), 1734–1766 (2015)
17. Emoto, K., Fischer, S., Hu, Z.: Generate, test, and aggregate: A calculation-based framework for systematic parallel programming with MapReduce. In: *ESOP*. pp. 254–273 (2012)
18. Herodotou, H., Babu, S.: Profiling, what-if analysis, and cost-based optimization of MapReduce programs. *Proceedings of the VLDB Endowment* 4(11), 1111–1122 (2011)
19. Herodotou, H., Borisov, N., Babu, S.: Query optimization techniques for partitioned tables. pp. 49–60. *SIGMOD '11*
20. Ioannidis, Y.E.: Query optimization. *ACM Comput. Surv.* 28(1), 121–123 (1996)
21. Karloff, H., Suri, S., Vassilvitskii, S.: A model of computation for MapReduce. In: *SODA*. pp. 938–948 (2010)
22. Leijen, D., Föhndrich, M., Burckhardt, S.: Prettier concurrency: Purely functional concurrent revisions. In: *Haskell*. pp. 83–94 (2011)
23. Liu, C., Zhang, J., Zhou, H., McDirmid, S., Guo, Z., Moscibroda, T.: Automating distributed partial aggregation. In: *SoCC*. pp. 1:1–1:12 (2014)
24. Radoi, C., Fink, S.J., Rabbah, R.M., Sridharan, M.: Translating imperative code to MapReduce. In: *OOPSLA*. pp. 909–927 (2014)
25. Sakr, S., Liu, A., Fayoumi, A.G.: The family of MapReduce and large-scale data processing systems. *ACM Comput. Surv.* 46(1), 11:1–11:44 (2013)
26. Tian, Y., Tatikonda, S., Reinwald, B.: Scalable and numerically stable descriptive statistics in SystemML. In: *ICDE*. pp. 1351–1359 (2012)
27. Xiao, T., Zhang, J., Zhou, H., Guo, Z., McDirmid, S., Lin, W., Chen, W., Zhou, L.: Non-determinism in MapReduce considered harmful? an empirical study on non-commutative aggregators in MapReduce programs. In: *Companion Proceedings of ICSE*. pp. 44–53 (2014)
28. Xu, Z., Hirzel, M., Rothermel, G.: Semantic characterization of MapReduce workloads. In: *IISWC*. pp. 87–97 (2013)
29. Zaharia, M., Chowdhury, M., Das, T., Dave, A., Ma, J., McCauly, M., Franklin, M.J., Shenker, S., Stoica, I.: Resilient distributed datasets: A fault-tolerant abstraction for in-memory cluster computing. In: *NSDI*. pp. 15–28 (2012)
30. Zaharia, M., Xin, R.S., Wendell, P., Das, T., Armbrust, M., Dave, A., Meng, X., Rosen, J., Venkataraman, S., Franklin, M.J., Ghodsi, A., Gonzalez, J., Shenker, S., Stoica, I.: Apache Spark: A unified engine for big data processing. *Commun. ACM* 59(11), 56–65 (Oct 2016)
31. Zhang, Z., Cherkasova, L., Verma, A., Loo, B.T.: Performance modeling and optimization of deadline-driven Pig programs. *ACM Trans. Auton. Adapt. Syst.* 8(3), 14:1–14:28 (2013)