# Equational reasoning for non-determinism monad:
# the case of Spark aggregation

**Shin-Cheng Mu**

# Equational Reasoning for Non-determinism Monad:

The Case of Spark Aggregation

SHIN-CHENG MU, Academia Sinica, Taiwan

As part of the author's studies on equational reasoning for monadic programs, this report focus on non-determinism monad. We discuss what properties this monad should satisfy, what additional operators and notations can be introduced to facilitate equational reasoning about non-determinism, and put them to the test by proving a number of properties in our example problem inspired by the author's previous work on proving properties of Spark aggregation.

## 1 INTRODUCTION

In functional programming, *pure* programs are those that can be understood as static mappings from inputs to outputs. The main advantage of staying in the pure realm is that properties of pure entities can be proved by equational reasoning. Side effects, in contrast, used to be considered the "awkward squad" that are difficult to be reasoned about. Gibbons and Hinze [2011], however, showed that effectful, monadic programs may also be reasoned about in a mathematical manner, using monad laws and properties of effect operators.

This report is part of a series of the author's studies on equational reasoning for monadic programs. In this report we focus on non-determinism monad — in our definition that is a monad having two effect operators, one allowing a program to fail, another allowing a non-deterministic choice between two results. We discuss what properties these operators should satisfy, what additional operators and notations can be introduced to facilitate equational reasoning of this monad, and put them to the test by proving a number of properties in our example problem: Spark aggregation.

Much of this report is inspired by the author's joint work with Chen et al. [2017], in which we formalised Spark, a platform for distributed computation, and derived properties under which a distributed Spark aggregation represents a *deterministic* computation. Therefore, many examples in this report are about finding out when processing a non-deterministic permutation (simulating arbitrary distribution of data) produces a deterministic result.

## 2 MONAD AND NON-DETERMINISM

A monad consists of a type constructor $M :: * \rightarrow *$ and two operators $return :: a \rightarrow M\ a$ and "bind" $(\lll) :: (a \rightarrow M\ b) \rightarrow M\ a \rightarrow M\ b$ that satisfy the following *monad laws*:

$$f \lll return\ x = f\ x\ , \tag{1}$$

$$return \lll m = m\ , \tag{2}$$

$$f \lll (g \lll m) = (\lambda x \rightarrow f \lll g\ x) \lll m\ . \tag{3}$$

Rather than the usual $(\ggg) :: M\ a \rightarrow (a \rightarrow M\ b) \rightarrow M\ b$, in the laws above we use the reversed bind $(\lll)$, which is consistent with the direction of function composition and more readable when we program in a style that uses composition. When we use bind with $\lambda$-abstractions, it is more natural to write $m \ggg \lambda x \rightarrow f\ x$. In this report we use the former more than the latter, thus the choice of notation. We also define $m_1 \ll m_2 = const\ m_1 \lll m_2$. Note that $(\gg)$ has type $M\ a \rightarrow M\ b \rightarrow M\ b$.

More operators we find useful are given in Figure 1. Right-to-left Kleisli composition, denoted by $(\lll)$, composes two monadic operations $a \rightarrow M\ b$ and $b \rightarrow M\ c$ into an operation $a \rightarrow M\ c$. Operators $(\$)$ and $(\bullet)$ are monadic counterparts of function application and composition: $(\$)$ applies a pure function to a monad, while $(\bullet)$ composes a pure function after a monadic function.

$$(\lll) :: (b \to M\ c) \to (a \to M\ b) \to a \to M\ c$$
$$(f \lll g)\ x = f \lll g\ x$$
$$(\$\!\!\$) :: (a \to b) \to M\ a \to M\ b$$
$$f \$\!\!\$\ m = (return \cdot f) \lll m$$
$$(\bullet) :: (b \to c) \to (a \to M\ b) \to (a \to M\ c)$$
$$f \bullet g = (return \cdot f) \lll g$$

Fig. 1. Some monadic operators we find handy for this paper.

We now introduce a collections of properties that allows us to rotate an expression that involves two operators and three operands. These properties will be handy when we need to move parenthesis around in expressions. To begin with, the following properties show that $(\$\!\!\$)$ and $(\bullet)$ share properties similar to pure function application and composition:

$$(f \bullet g)\ x = f \$\!\!\$\ g\ x\,, \tag{4}$$

$$f \$\!\!\$\ (g \$\!\!\$\ m) = (f \cdot g) \$\!\!\$\ m\,, \tag{5}$$

$$f \bullet (g \bullet h) = (f \cdot g) \bullet h\,. \tag{6}$$

We also have the following law that allows us to rotate an expression that uses $(\bullet)$ and $(\cdot)$:

$$f \bullet (g \cdot h) = (f \bullet g) \cdot h\,. \tag{7}$$

Note that $g$ in (7) must be a function returning a monad. Furthermore, (8) and (9) relate $(\lll)$ and $(\$\!\!\$)$, both operators applying functions to monads, while (10) and (11) relate $(\lll)$ and $(\bullet)$, both operators composing functions on monads:

$$f \lll (g \$\!\!\$\ m) = (f \cdot g) \lll m\,, \tag{8}$$

$$f \$\!\!\$\ (g \lll m) = (f \bullet g) \lll m\,, \tag{9}$$

$$f \lll (g \bullet h) = (f \cdot g) \lll h\,, \tag{10}$$

$$f \bullet (g \lll h) = (f \bullet g) \lll h\,. \tag{11}$$

Having these properties is one of the advantages of writing $(\lll)$ and $(\lll)$ backwards. All the properties above can be proved by expanding definitions, and it is a good warming-up exercise proving some of them. Some of them are proved in Appendix A.

None of these operators and properties are strictly necessary: they can all be reduced to *return*, $(\lll)$, and $\lambda$-abstractions. As is often the case when designing notations, having more operators allows ideas to be expressed concisely in a higher level of abstraction, at the expense of having more properties to memorise. It is personal preference where the balance should be. Properties (4) through (11) may look like a lot of properties to remember. In practice, we find it usually sufficient to let us be guided by types. For example, when we have $f \$\!\!\$\ g\ x$ and want to bring $f$ and $g$ together, by their types we can figure out the resulting expression should be $(f \bullet g)\ x$.

*Non-determinism Monad.* Non-determinism is the only effect we use in this report. We assume two operators $\emptyset$ and $(\|)$: the former denotes failure, while $m \parallel n$ denotes that the computation may yield either $m$ or $n$. As pointed out by Gibbons and Hinze [2011], for proofs and derivations, what matters is not how a monad is implemented but what properties its operators satisfy. What laws $\emptyset$ and $(\|)$ should satisfy, however, can be a tricky issue. As discussed by Kiselyov [2015], it eventually comes down to what we use the monad for. It is usually expected that $(a, (\|), \emptyset)$ be a monoid. That

is, ($[\!]$) is associative, with $\emptyset$ as its zero:

$$(m \mathbin{[\!]} n) \mathbin{[\!]} k \;=\; m \mathbin{[\!]} (n \mathbin{[\!]} k)\,,$$

$$\emptyset \mathbin{[\!]} m \;=\; m \;=\; m \mathbin{[\!]} \emptyset\,.$$

It is also assumed that monadic bind distributes into ($[\!]$) from the end, while $\emptyset$ is a right zero for ($\gg\!\!\ll$):

$$f \gg\!\!\ll (m_1 \mathbin{[\!]} m_2) \;=\; (f \gg\!\!\ll m_1) \mathbin{[\!]} (f \gg\!\!\ll m_2)\,, \tag{12}$$

$$f \gg\!\!\ll \emptyset \;=\; \emptyset\,. \tag{13}$$

For our purpose in this section, we also assume that ($[\!]$) is commutative ($m \mathbin{[\!]} n = n \mathbin{[\!]} m$) and idempotent ($m \mathbin{[\!]} m = m$). Implementation of such non-determinism monads have been studied by Fischer et al. [2011].

Here are some induced laws about how ($\$$) interacts with *return* and non-determinism operators:

$$f \mathbin{\$} return\ x = return\ (f\ x)\,, \tag{14}$$

$$f \mathbin{\$} \emptyset = \emptyset\,, \tag{15}$$

$$f \mathbin{\$} (m_1 \mathbin{[\!]} m_2) = (f \mathbin{\$} m_1) \mathbin{[\!]} (f \mathbin{\$} m_2)\,. \tag{16}$$

## 3 PERMUTATION AND INSERTION

As a warm-up example, the function *perm* non-deterministically computes a permutation of its input, using an auxiliary function *insert* that inserts an element to an arbitrary position in a list:

$$
\begin{aligned}
&perm && :: [a] \rightarrow \mathsf{M}\ [a] \\
&perm\ [\,] && = return\ [\,] \\
&perm\ (x:xs) && = insert\ x \gg\!\!\ll perm\ xs\,, \\[4pt]
&insert && :: a \rightarrow [a] \rightarrow \mathsf{M}\ [a] \\
&insert\ x\ [\,] && = return\ [x] \\
&insert\ x\ (y:xs) && = return\ (x:y:xs) \mathbin{[\!]} ((y:) \mathbin{\$} insert\ x\ xs)\,.
\end{aligned}
$$

For example, possible results of *perm* $[0, 1, 2]$ include $[0, 1, 2]$, $[0, 2, 1]$, $[1, 0, 2]$, $[1, 2, 0]$, $[2, 0, 1]$, and $[2, 1, 0]$.

*Determinism.* The following lemma presents properties under which permuting the input list does not change the result of a *foldr*:

LEMMA 3.1. *Given* $(\odot) :: a \rightarrow b \rightarrow b$. *If* $x \odot (y \odot z) = y \odot (x \odot z)$ *for all* $x, y :: a$ *and* $z :: b$, *we have*

$$foldr\ (\odot)\ z \mathbin{\langle\!\bullet\rangle} perm = return \cdot foldr\ (\odot)\ z\,.$$

Since *perm* is defined in terms of *insert*, proof of Lemma 3.1 naturally depends on a lemma about a related property of *insert*:

LEMMA 3.2. *Given* $(\odot) :: a \rightarrow b \rightarrow b$, *we have*

$$foldr\ (\odot)\ z \mathbin{\langle\!\bullet\rangle} insert\ x = return \cdot foldr\ (\odot)\ z \cdot (x:)\,,$$

*provided that* $x \odot (y \odot z) = y \odot (x \odot z)$ *for all* $x, y :: a$ *and* $z :: b$.

PROOF. Prove $foldr\ (\odot)\ z \mathbin{\$} insert\ x\ xs = return\ (foldr\ (\odot)\ z\ (x:xs))$. Induction on *xs*.
CASE $xs := [\,]$:

$$
\begin{aligned}
&\quad foldr\ (\odot)\ z \mathbin{\$} insert\ x\ [\,] \\
&= \quad \{\text{definition of } insert\ \}
\end{aligned}
$$

  $foldr\;(\odot)\;z\;⟨\$⟩\;return\;[x]$
$=$ { by (14) }
  $return\;(foldr\;(\odot)\;z\;[x])\;.$

Case $xs := y : xs$:

  $foldr\;(\odot)\;z\;⟨\$⟩\;insert\;x\;(y : xs)$
$=$ { definition of $insert$ }
  $foldr\;(\odot)\;z\;⟨\$⟩\;(return\;(x : y : xs)\;[\!]\;((y{:})\;⟨\$⟩\;insert\;x\;xs))$
$=$ { by (16), (14), and (5) }
  $return\;(foldr\;(\odot)\;z\;(x : y : xs))\;[\!]\;((foldr\;(\odot)\;z \cdot (y{:}))\;⟨\$⟩\;insert\;x\;xs)\;.$

Focus on the second branch of ($[\!]$):

  $(foldr\;(\odot)\;z \cdot (y{:}))\;⟨\$⟩\;insert\;x\;xs$
$=$ { definition of $foldr$ }
  $((y\odot) \cdot foldr\;(\odot)\;z)\;⟨\$⟩\;insert\;x\;xs$
$=$ { by (5) }
  $(y\odot)\;⟨\$⟩\;(foldr\;(\odot)\;z\;⟨\$⟩\;insert\;x\;xs)$
$=$ { induction }
  $(y\odot)\;⟨\$⟩\;return\;(foldr\;(\odot)\;z\;(x : xs))$
$=$ { by (14) }
  $return\;(y \odot foldr\;(\odot)\;z\;(x : xs))$
$=$ { definition of $foldr$ }
  $return\;(y \odot (x \odot foldr\;(\odot)\;z\;xs))$
$=$ { since $x \odot (y \odot z) = y \odot (x \odot z)$ }
  $return\;(foldr\;(\odot)\;z\;(x : y : xs))\;.$

Thus we have

  $(foldr\;(\odot)\;z\;⟨\bullet⟩\;insert\;x)\;(y : xs)$
$=$ { calculation above }
  $return\;(foldr\;(\odot)\;z\;(x : y : xs))\;[\!]\;return\;(foldr\;(\odot)\;z\;(x : y : xs))$
$=$ { idempotence of ($[\!]$) }
  $return\;(foldr\;(\odot)\;z\;(x : y : xs))\;.$

                          $\square$

Proof of Lemma 3.1 then follows:

Proof. Prove that $foldr\;(\odot)\;z\;⟨\$⟩\;perm\;xs = return\;(foldr\;(\odot)\;z\;xs)$. Induction on $xs$.
Case $xs := [\,]$:

  $foldr\;(\odot)\;z\;⟨\$⟩\;perm\;[\,]$
$=$ { definitions of $perm$ }
  $foldr\;(\odot)\;z\;⟨\$⟩\;return\;[\,]$
$=$ { by (14) }
  $return\;(foldr\;(\odot)\;z\;[\,])\;.$

Case $xs := x : xs$:

  $foldr\;(\odot)\;z\;⟨\$⟩\;perm\;(x : xs)$
$=$ { definition of $perm$ }

$$foldr \ (\odot) \ z \ \langle\$\rangle \ (insert \ x \gg\!\!\!= perm \ xs)$$
$$= \quad \{ \text{by (9)} \}$$
$$(foldr \ (\odot) \ z \ \langle\bullet\rangle \ insert \ x) \gg\!\!\!= perm \ xs$$
$$= \quad \{ \text{Lemma 3.2} \}$$
$$(return \cdot foldr \ (\odot) \ z \cdot (x:)) \gg\!\!\!= perm \ xs$$
$$= \quad \{ \text{definitions of } foldr \text{ and } (\langle\$\rangle) \}$$
$$((x\odot) \cdot foldr \ (\odot) \ z) \ \langle\$\rangle \ perm \ xs$$
$$= \quad \{ \text{by (5)} \}$$
$$(x\odot) \ \langle\$\rangle \ (foldr \ (\odot) \ z \ \langle\$\rangle \ perm \ xs)$$
$$= \quad \{ \text{induction} \}$$
$$(x\odot) \ \langle\$\rangle \ (return \ (foldr \ (\odot) \ z \ xs))$$
$$= \quad \{ \text{by (14)} \}$$
$$return \ (x \odot foldr \ (\odot) \ z \ xs)$$
$$= \quad \{ \text{definition of } foldr \}$$
$$return \ (foldr \ (\odot) \ z \ (x : xs)) \ .$$

$\square$

*Map, Filter, and Permutation.* It is not hard for one to formulate the following relationship between *map* and *perm*, which is also based on a related property relating *map* and *insert*:[1]

LEMMA 3.3. *perm* · *map f* = *map f* ⟨•⟩ *perm*.

LEMMA 3.4. *insert* (*f x*) · *map f* = *map f* ⟨•⟩ *insert x*.

The lemma is true because *map f* is a pure computation — in reasoning about monadic programs it is helpful, and sometimes essential, to identify its pure segments, because these are the parts more properties are applicable. Note that the composition (·) on the lefthand side is turned into (⟨•⟩) once we move *map f* leftwards.

We prove only Lemma 3.4.

PROOF. Prove by induction on *xs* that *map f* ⟨\$⟩ *insert x xs* = *insert* (*f x*) (*map f xs*) for all *xs*. We present only the inductive case *xs* := *y* : *xs*:

$$map \ f \ \langle\$\rangle \ insert \ x \ (y : xs)$$
$$= \quad \{ \text{definition of } insert \}$$
$$map \ f \ \langle\$\rangle \ (return \ (x : y : xs) \ [\!] \ ((y:) \ \langle\$\rangle \ insert \ x \ xs))$$
$$= \quad \{ \text{by (16) and (14)} \}$$
$$return \ (map \ f \ (x : y : xs)) \ [\!] \ (map \ f \ \langle\$\rangle \ ((y:) \ \langle\$\rangle \ insert \ x \ xs)) \ .$$

For the second branch we reason:

$$map \ f \ \langle\$\rangle \ ((y:) \ \langle\$\rangle \ insert \ x \ xs)$$
$$= \quad \{ \text{by (5)} \}$$
$$(map \ f \cdot (y:)) \ \langle\$\rangle \ insert \ x \ xs$$
$$= \quad \{ \text{definition of } map \}$$
$$((f \ y:) \cdot map \ f) \ \langle\$\rangle \ insert \ x \ xs$$
$$= \quad \{ \text{by (5)} \}$$
$$(f \ y:) \ \langle\$\rangle \ (map \ f \ \langle\$\rangle \ insert \ x \ xs)$$

---

[1]Lemma 3.3 and 3.4 are in fact free theorems of *perm* and *insert* [Voigtländer 2009]. They serve as good exercises, nevertheless.

$=$    { induction }
  $(f\ y\!:) \circledS insert\ (f\ x)\ (map\ f\ xs)$ .

Thus we have

  $map\ f \circledS insert\ x\ (y:xs)$
$=$    { calculation above }
  $return\ (f\ x:f\ y:map\ f\ xs) \llbracket ((f\ y\!:) \circledS (insert\ (f\ x)\ (map\ f\ xs)))$
$=$    { definitions of *insert* and *map* }
  $insert\ (f\ x)\ (map\ f\ (y:xs))$ .

$\square$

One may have noticed that the style of proof is familiar: replace *return x* by $[x]$ and $(\llbracket)$ by $(+\!+)$, the proof is more-or-less what one would do for a list version of *insert*. This is exactly the point: the style of proofs we use to do for pure programs still works for monadic programs, as long as the monad satisfies the demanded laws, be it a list, a more advanced implementation of non-determinism, or a monad having other effects.

A similar property relating *perm* and *filter* can be formulated.

LEMMA 3.5.  $perm \cdot filter\ p = filter\ p \diamond perm.$

Its proof is routine and omitted. Finally, in a number of occasions it helps to know that *xs* is a result of *perm xs*. The proof is also routine and omitted.

LEMMA 3.6.  *For all xs we have that* $perm\ xs = return\ xs \llbracket m$ *for some m.*

## 4   SPARK AGGREGATION

Spark [Zaharia et al. 2012] is a popular platform for scalable distributed data-parallel computation based on a flexible programming environment with high-level APIs, considered by many as the successor of MapReduce. In a typical Spark program, data is partitioned and stored distributively on read-only *Resilient Distributed Datasets* (RDDs) — we can think of it as a list of lists, where each sub-list is potentially stored on a remote node. On an RDD one can apply operations, called *combinators*, such as *map*, *reduce*, and *aggregate*. The *aggregate* combinator, for example, takes user-defined functions $(\otimes)$ and $(\oplus)$: $(\otimes)$ accumulates a sub-result for each data partition while $(\oplus)$ merges sub-results across different partitions.

Programming in Spark, however, can be tricky. Since sub-results are computed across partitions concurrently, the order of their applications varies on different executions. Aggregation in Spark is therefore inherently non-deterministic. An example from Chen et al. [2017] showed that computing the integral of $x^{73}$ from $x = -2$ to $x = 2$, which should be $0$, using a function in the Spark machine learning library, yields results ranging from $-8192.0$ to $12288.0$ in different runs. It is thus desirable to find out conditions, which Spark's documentation does not specify formally, under which a Spark computation yields deterministic outcomes.

### 4.1   List Homomorphism

Since a Spark aggregation is typically used to computes a *list homomorphism* [Bird 1987], we digress a little in this section to give a brief review and present some results that we will use. A function $h :: List\ a \rightarrow b$ is called a list homomorphism if there exists $z :: b$, $k :: a \rightarrow b$, and $(\oplus) :: b \rightarrow b \rightarrow b$ such that:

$$h\ [\ ]\qquad = z$$
$$h\ [x]\qquad = k\ x$$
$$h\ (xs \mathbin{+\!\!+} ys) = h\ xs \oplus h\ ys\ .$$

That $h$ is such a list homomorphism is denoted by $h = hom\ (\oplus)\ k\ z$. Note that the properties above implicitly demand that $(\oplus)$ be associative with $z$ as its identity element.

Lemma 4.1 and 4.2 below are about when a computation defined in terms of *foldr* is actually a list homomorphism. In Lemma 4.2, *img f* denotes the image of a function $f$.

LEMMA 4.1. $h = hom\ (\oplus)\ (h \cdot wrap)\ z$ *if and only if* $foldr\ (\oplus)\ z \cdot map\ h = h \cdot concat$, *where* $wrap\ x = [x]$.

LEMMA 4.2. *Let* $(\oplus) :: b \to b \to b$ *be associative on* $img\ (foldr\ (\otimes)\ z)$ *with $z$ as its identity, where* $(\otimes) :: a \to b \to b$. *We have* $foldr\ (\otimes)\ z = hom\ (\oplus)\ (\otimes z)\ z$ *if and only if* $x \otimes (y \oplus w) = (x \otimes y) \oplus w$ *for all* $x :: a$ *and* $y, w \in img\ (foldr\ (\otimes)\ z)$.

Notice, in Lemma 4.2, that $(\otimes z) = foldr\ (\otimes)\ z \cdot wrap$. Proofs of both lemmas are interesting exercises, albeit being a bit off-topic. They are recorded in Appendix A.

### 4.2 Formalisation and Results

Distributed collections of data are represented by *Resilient Distributed Datasets* (RDDs) in Spark. Informally, an RDD is a collection of data entries; these data entries are further divided into partitions stored on different machines. Abstractly, an RDD can be seen as a list of lists:

> **type** Partition $a$ = $[a]$ ,
> **type** RDD $a$    = $[$Partition $a]$ ,

where each Partition may be stored in a different machine.

While Spark provides a collection of *combinators* (functions on RDDs that are designed to be composed to form larger programs), in this report we focus on a particular one, *aggregate*. It can be seen as a parallel implementation *foldr*. The combinator processes an RDD in two levels: each partition is first processed locally on one machine by $foldr\ (\otimes)\ z$. The sub-results are then communicated and combined — this second step can be think of as another *foldr* with $(\oplus)$.[2]

Spark programmers like to assume that their programs are deterministic. To exploit concurrency, however, the sub-results from each machine might be processed in arbitrary order and the result could be non-deterministic. The following is our characterisation of *aggregate*, where we use *perm* to model the fact that sub-results from each machine are processed in unknown order:

> $aggregate :: b \to (a \to b \to b) \to (b \to b \to b) \to RDD\ a \to M\ b$
> $aggregate\ z\ (\otimes)\ (\oplus) = foldr\ (\oplus)\ z \mathbin{\diamond} (perm \cdot map\ (foldr\ (\otimes)\ z))\ .$

It is clear from the types that $foldr\ (\otimes)\ z$ and $foldr\ (\oplus)\ z$ are pure computations, and non-determinism is introduced solely by *perm*.

*Deterministic Aggregation.* We are interested in finding out conditions under which *aggregate* produces deterministic outcomes.

THEOREM 4.3. *Given* $(\otimes) :: a \to b \to b$ *and* $(\oplus) :: b \to b \to b$, *where* $(\oplus)$ *is associative and commutative, we have:*

> $aggregate\ z\ (\otimes)\ (\oplus) = return \cdot foldr\ (\oplus)\ z \cdot map\ (foldr\ (\otimes)\ z)\ .$

---

[2]In fact, the actual Spark aggregation (and that modelled in Chen et al. [2017]) are like *foldl*. For convenience in our proofs we see all list operations the other way round and use *foldr*. This is not a fundamental difference.

Proof. We reason:

$aggregate\ z\ (\otimes)\ (\oplus)$

$=$    { definition of $aggregate$ }

$foldr\ (\oplus)\ z\ \diamond\ (perm \cdot map\ (foldr\ (\otimes)\ z))$

$=$    { by (7) }

$(foldr\ (\oplus)\ z\ \diamond\ perm) \cdot map\ (foldr\ (\otimes)\ z)$

$=$    { Lemma 3.1, since $(\oplus)$ is associative and commutative }

$return \cdot foldr\ (\oplus)\ z \cdot map\ (foldr\ (\otimes)\ z)$ .

□

The following corollary summaries the results and present conditions under which *aggregate* computes a homomorphism.

COROLLARY 4.4. *$aggregate\ z\ (\otimes)\ (\oplus) = return \cdot hom\ (\oplus)\ (\otimes z)\ z \cdot concat$, provided that $(\oplus)$ is associative, commutative, and has $z$ as identity, and that $x \otimes (y \oplus w) = (x \otimes y) \oplus w$ for all $x :: a$ and $y, w \in img\ (foldr\ (\otimes)\ z)$.*

Proof. We reason:

$aggregate\ z\ (\otimes)\ (\oplus)$

$=$    { Theorem 4.3 }

$return \cdot foldr\ (\oplus)\ z \cdot map\ (foldr\ (\otimes)\ z)$

$=$    { $foldr\ (\otimes)\ z = hom\ (\oplus)\ (\otimes z)\ z$ by Lemma 4.2; Lemma 4.1 }

$return \cdot hom\ (\oplus)\ (\otimes z)\ z \cdot concat$ .

□

*Determinism Implies Homomorphism.* The final part of the report deals with an opposite question: what can we infer if we know that *aggregate* is deterministic? To answer that, however, we need to assume two more properties:

$$m_1 \ [\!]\ m_2 = return\ x \implies m_1 = m_2 = return\ x. \tag{17}$$

$$return\ x_1 = return\ x_2 \implies x_1 = x_2. \tag{18}$$

Property (17) can be seen as the other direction of idempotency of $([\!])$, while (18) states that *return* is injective.

The following lemma can be understood this way: when $aggregate\ z\ (\otimes)\ (\oplus)$, which could be non-deterministic, can be performed by a deterministic function, the operator $(\oplus)$ should be insensitive to ordering:

LEMMA 4.5. *If $aggregate\ z\ (\otimes)\ (\oplus) = return \cdot foldr\ (\otimes)\ z \cdot concat$, and $perm\ xss = return\ yss\ [\!]\ m$ for some $m$, we have*

$foldr\ (\otimes)\ z\ (concat\ xss) =$

$\quad foldr\ (\oplus)\ z\ (map\ (foldr\ (\otimes)\ z)\ xss) =$

$\quad\quad foldr\ (\oplus)\ z\ (map\ (foldr\ (\otimes)\ z)\ yss)$ .

Proof. We reason:

$return \cdot foldr\ (\otimes)\ z \cdot concat\ \$\ xss$

$=$    { assumption }

$aggregate\ z\ (\otimes)\ (\oplus)\ \$\ xss$

$=$   { definition of *aggregate*, Lemma 3.3, and (6) }
   $(foldr\ (\oplus)\ z \cdot map\ (foldr\ (\otimes)\ z)) \mathrel{\langle\$\rangle} perm\ xss$
$=$   { assumption: $perm\ xss = return\ yss \mathbin{[\!]} m$, by (16) and (14) }
   $(return \cdot foldr\ (\oplus)\ z \cdot map\ (foldr\ (\otimes)\ z) \mathbin{\$} yss) \mathbin{[\!]}$
      $((foldr\ (\oplus)\ z \cdot map\ (foldr\ (\otimes)\ z)) \mathrel{\langle\$\rangle} m)\ .$

Thus by (17) and (18), $foldr\ (\otimes)\ z \cdot concat \mathbin{\$} xss$ equals $foldr\ (\oplus)\ z \cdot map\ (foldr\ (\otimes)\ z) \mathbin{\$} yss$. The former also equals $foldr\ (\oplus)\ z \cdot map\ (foldr\ (\otimes)\ z) \mathbin{\$} xss$ because, by Lemma 3.6, $perm\ xss = return\ xss \mathbin{[\!]} m$ for some $m$.                                                                                                   □

Based on Lemma 4.5, the following theorem explicitly states that $(\oplus)$ should be associative, commutative, and has $z$ as its identity in restricted domain.

THEOREM 4.6. *If aggregate* $z\ (\otimes)\ (\oplus) = return \cdot foldr\ (\otimes)\ z \cdot concat$, *we have that* $(\oplus)$, *when restricted to values in* $img\ (foldr\ (\otimes)\ z)$, *is associative, commutative, and has* $z$ *as its identity.*

PROOF. In the discussion below, let $x, y$, and $w$ be in $img\ (foldr\ (\otimes)\ z)$. That is, there exists $xs, ys$, and $ws$ such that $x = foldr\ (\otimes)\ z\ xs$, $y = foldr\ (\otimes)\ z\ ys$, and $w = foldr\ (\otimes)\ z\ ws$.
IDENTITY. We reason:

   $y$
$= foldr\ (\otimes)\ z\ (concat\ [xs])$
$=$   { $perm\ [xs] = return\ [xs] \mathbin{[\!]} \emptyset$, Lemma 4.5 }
   $foldr\ (\oplus)\ z\ (map\ (foldr\ (\otimes)\ z)\ [xs])$
$= y \oplus z\ .$

Thus $z$ is a right identity of $(\oplus)$. Similarly,

   $y$
$= foldr\ (\otimes)\ z\ (concat\ [[\,], xs])$
$=$   { $perm\ [[\,], xs] = return\ [[\,], xs] \mathbin{[\!]} m$, Lemma 4.5 }
   $foldr\ (\oplus)\ z\ (map\ (foldr\ (\otimes)\ z)\ [[\,], xs])$
$= z \oplus (y \oplus z)$
$=$   { $z$ is a right identity of $(\oplus)$ }
   $z \oplus y\ .$

Thus $z$ is also a left identity of $(\oplus)$.
COMMUTATIVITY. We reason:

   $x \oplus y$
$=$   { $z$ is a right identity }
   $x \oplus (y \oplus z)$
$= foldr\ (\oplus)\ z\ (map\ (foldr\ (\otimes)\ z)\ [xs, ys])$
$=$   { $perm\ [xs, ys] = return\ [ys, xs] \mathbin{[\!]} m$, Lemma 4.5 }
   $foldr\ (\oplus)\ z\ (map\ (foldr\ (\otimes)\ z)\ [ys, xs])$
$= y \oplus (x \oplus z)$
$=$   { $z$ is a right identity }
   $y \oplus x\ .$

ASSOCIATIVITY. We reason:

   $x \oplus (y \oplus w)$
$=$   { $z$ is a right identity }

$$x \oplus (y \oplus (w \oplus z))$$
$$= foldr\ (\oplus)\ z\ (map\ (foldr\ (\otimes)\ z)\ [\,xs, ys, ws\,])$$
$$= \quad \{\ (\oplus)\ \text{commutative}\ \}$$
$$foldr\ (\oplus)\ z\ (map\ (foldr\ (\otimes)\ z)\ [\,ws, xs, ys\,])$$
$$= w \oplus (x \oplus (y \oplus z))$$
$$= \quad \{\ z\ \text{is a right identity}\ \}$$
$$w \oplus (x \oplus y)$$
$$= \quad \{\ (\oplus)\ \text{commutative}\ \}$$
$$(x \oplus y) \oplus w\ .$$

$\square$

**Theorem 4.7.** *If* $aggregate\ z\ (\otimes)\ (\oplus) = return \cdot foldr\ (\otimes)\ z \cdot concat$, *we have* $foldr\ (\otimes)\ z = hom\ (\oplus)\ (\otimes z)\ z$.

**Proof.** Apparently $foldr\ (\otimes)\ z\ [\,]\ = z$ and $foldr\ (\otimes)\ z\ [\,x\,] = x \otimes z$. We are left with proving the case for concatenation.

$$foldr\ (\otimes)\ z\ (xs \mathbin{+\!\!+} ys)$$
$$= foldr\ (\otimes)\ z\ (concat\ [\,xs, ys\,])$$
$$= \quad \{\ \text{Lemma 4.5}\ \}$$
$$foldr\ (\oplus)\ z\ (map\ (foldr\ (\otimes)\ z)\ [\,xs, ys\,])$$
$$= foldr\ (\otimes)\ z\ xs \oplus (foldr\ (\otimes)\ z\ ys \oplus z)$$
$$= \quad \{\ \text{Theorem 4.6},\ z\ \text{is identity}\ \}$$
$$foldr\ (\otimes)\ z\ xs \oplus foldr\ (\otimes)\ z\ ys\ .$$

$\square$

**Corollary 4.8.** *Given* $(\otimes) :: a \rightarrow b \rightarrow b$ *and* $(\oplus) :: b \rightarrow b \rightarrow b$. $aggregate\ z\ (\otimes)\ (\oplus) = return \cdot foldr\ (\otimes)\ z \cdot concat$ *if and only if* $(img\ (foldr\ (\otimes)\ z), (\oplus), z)$ *forms a commutative monoid, and that* $foldr\ (\otimes)\ z = hom\ (\oplus)\ (\otimes z)\ z$.

**Proof.** A conclusion following from Theorem 4.3, Theorem 4.6, and Theorem 4.7. $\square$

## REFERENCES

Reynald Affeldt, David Nowak, and Takafumi Saikawa. 2019. A hierarchy of monadic effects for program verification using equational reasoning. In *Mathematics of Program Construction*, Graham Hutton (Ed.). Springer.

Richard S. Bird. 1987. An introduction to the theory of lists. In *Logic of Programming and Calculi of Discrete Design*, Manfred Broy (Ed.). Number 36 in NATO ASI Series F. Springer-Verlag, 3–42.

Yu-Fang Chen, Chih-Duo Hong, Ondřej Lengál, Shin-Cheng Mu, Nishant Sinha, and Bow-Yaw Wang. 2017. An executable sequential specification for Spark aggregation. In *International Conference on Networked Systems*. Springer-Verlag.

Sebastian Fischer, Oleg Kiselyov, and Chung-chieh Shan. 2011. Purely functional lazy nondeterministic programming. *Journal of Functional Programming* 21, 4-5 (September 2011), 413–465.

Jeremy Gibbons and Ralf Hinze. 2011. Just do it: simple monadic equational reasoning. In *International Conference on Functional Programming*, Olivier Danvy (Ed.). ACM Press, 2–14.

Oleg Kiselyov. 2015. Laws of MonadPlus. http://okmij.org/ftp/Computation/monads.html#monadplus.

Janis Voigtländer. 2009. Free theorems involving type constructor classes. In *International Conference on Functional Programming*, Andrew Tolmach (Ed.). ACM Press, 173–184.

Matei Zaharia, Mosharaf Chowdhury, Tathagata Das, Ankur Dave, Justin Ma, Murphy McCauley, Michael J. Franklin, Scott Shenker, and Ion Stoica. 2012. Resilient distributed datasets: a fault-tolerant abstraction for in-memory cluster computing. In *Networked Systems Design and Implementation*, Steven Gribble and Dina Katabi (Eds.). USENIX.

## A  MISCELLANEOUS PROOFS

*Proving* (8). $f \lll (g \Diamond\$\Diamond m) = (f \cdot g) \lll m$.

PROOF. We reason:

$$
\begin{aligned}
& f \lll (g \Diamond\$\Diamond m) \\
= \quad & \{ \text{definition of } (\Diamond\$\Diamond) \} \\
& f \lll ((return \cdot g) \lll m) \\
= \quad & \{ \text{monad law (3)} \} \\
& (\lambda x \rightarrow f \lll return (g\ x)) \lll m \\
= \quad & \{ \text{monad law (1)} \} \\
& (\lambda x \rightarrow f (g\ x)) \lll m \\
= \ & (f \cdot g) \lll m\ .
\end{aligned}
$$

□

*Proving* (5). $f \Diamond\$\Diamond (g \Diamond\$\Diamond m) = (f \cdot g) \Diamond\$\Diamond m$.

PROOF. We reason:

$$
\begin{aligned}
& f \Diamond\$\Diamond (g \Diamond\$\Diamond m) \\
= \quad & \{ \text{definition of } (\Diamond\$\Diamond) \} \\
& (return \cdot f) \lll (g \Diamond\$\Diamond m) \\
= \quad & \{ \text{by (8)} \} \\
& (return \cdot f \cdot g) \lll m \\
= \quad & \{ \text{definition of } (\Diamond\$\Diamond) \} \\
& (f \cdot g) \Diamond\$\Diamond m\ .
\end{aligned}
$$

□

For the next results we prove a lemma:

$$
(f\lll) \cdot (g\lll) = (((f\lll) \cdot g)\lll)\ . \tag{19}
$$

$$
\begin{aligned}
& (f\lll) \cdot (g\lll) \\
= \quad & \{ \eta \text{ intro.} \} \\
& (\lambda m \rightarrow f \lll (g \lll m)) \\
= \quad & \{ \text{monad law (3)} \} \\
& (\lambda m \rightarrow (\lambda y \rightarrow f \lll g\ y) \lll m) \\
= \quad & \{ \eta \text{ reduction} \} \\
& (((f\lll) \cdot g)\lll)\ .
\end{aligned}
$$

*Proving* (6). $f \odot (g \odot m) = (f \cdot g) \odot m$.

PROOF. We reason:

$\quad f \odot (g \odot m)$
$= \quad \{ \text{definition of } (\odot) \}$
$\quad ((return \cdot f) \lll) \cdot ((return \cdot g) \lll) \cdot m$
$= \quad \{ \text{by (19)} \}$
$\quad ((((return \cdot f) \lll) \cdot return \cdot g) \lll) \cdot m$
$= \quad \{ \text{monad law (1)} \}$
$\quad ((return \cdot f \cdot g) \lll) \cdot m$
$= \quad \{ \text{definition of } (\odot) \}$
$\quad (f \cdot g) \odot m \ .$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \square$

*Proving* (10). $f \lll (g \odot h) = (f \cdot g) \lll h$.

PROOF. We reason:

$\quad f \lll (g \odot h)$
$= \quad \{ \text{definitions of } (\lll) \}$
$\quad (f \lll) \cdot ((return \cdot g) \lll) \cdot h$
$= \quad \{ \text{by (19)} \}$
$\quad (((f \lll) \cdot return \cdot g) \lll) \cdot h$
$= \quad \{ \text{monad law (1)} \}$
$\quad ((f \cdot g) \lll) \cdot h$
$= \quad \{ \text{definition of } (\lll) \}$
$\quad (f \cdot g) \lll h \ .$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \square$

*Proving* (11). $f \odot (g \lll h) = (f \odot g) \lll h$.

PROOF. We reason:

$\quad f \odot (g \lll h)$
$= \quad \{ \text{definitions of } (\lll) \text{ and } (\odot) \}$
$\quad ((return \cdot f) \lll) \cdot (g \lll) \cdot h$
$= \quad \{ \text{by (19)} \}$
$\quad ((((return \cdot f) \lll) \cdot g) \lll) \cdot h$
$= \quad \{ \text{definition of } (\odot) \}$
$\quad ((f \odot g) \lll) \cdot h$
$= \quad \{ \text{definition of } (\lll) \}$
$\quad (f \odot g) \lll h \ .$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \square$

*Proof of Lemma 4.1.*

PROOF. A Ping-pong proof.
DIRECTION ($\Rightarrow$). Let $h = hom\ (\oplus)\ (h \cdot wrap)\ z$, prove $foldr\ (\oplus)\ z\ (map\ h\ xss) = h\ (concat\ xss)$ by induction on $xss$.
CASE $xss := [\,]$:

$$foldr\ (\oplus)\ z\ (map\ h\ [\,])$$
$$= foldr\ (\oplus)\ z\ [\,]$$
$$= z$$
$$= h\ (concat\ [\,])\ .$$

CASE $xss := xs : xss$:

$$foldr\ (\oplus)\ z\ (map\ h\ (xs : xss))$$
$$= h\ xs \oplus foldr\ (\oplus)\ z\ (map\ h\ xss)$$
$$= \quad \{\,induction\,\}$$
$$h\ xs \oplus h\ (concat\ xss)$$
$$= \quad \{\,h\ homomorphism\,\}$$
$$h\ (concat\ (xs : xss))\ .$$

DIRECTION ($\Leftarrow$). Assuming $foldr\ (\oplus)\ z\ (map\ h\ xss) = h\ (concat\ xss)$, prove that $h = hom\ (\oplus)\ (h \cdot wrap)\ z$.
CASE empty list:

$$h\ [\,]$$
$$= h\ (concat\ [\,])$$
$$= \quad \{\,assumption\,\}$$
$$foldr\ (\oplus)\ z\ (map\ h\ [\,])$$
$$= z\ .$$

CASE singleton list: certainly $h\ [\,x\,] = h\ [\,x\,]$.
CASE concatentation:

$$h\ (xs + ys)$$
$$= h\ (concat\ [\,xs, ys\,])$$
$$= \quad \{\,assumption\,\}$$
$$foldr\ (\oplus)\ z\ (map\ h\ [\,xs, ys\,])$$
$$= h\ xs \oplus (h\ ys \oplus z)$$
$$= h\ xs \oplus h\ ys\ .$$

□

*Proof of Lemma 4.2.*

PROOF. A Ping-pong proof.
DIRECTION ($\Leftarrow$). We show that $foldr\ (\otimes)\ z = hom\ (\oplus)\ (\otimes z)\ z$, provided that $x \otimes (y \oplus w) = (x \otimes y) \oplus w$.
It is immediate that $foldr\ (\otimes)\ z\ [\,] = z$ around $foldr\ (\otimes)\ z\ [\,x\,] = x \otimes z$. For $xs + ys$, note that

$$foldr\ (\otimes)\ z\ (xs + ys) = foldr\ (\otimes)\ (foldr\ (\otimes)\ z\ ys)\ xs\ .$$

The aim is thus to prove that

$$foldr\ (\otimes)\ (foldr\ (\otimes)\ z\ ys)\ xs = (foldr\ (\otimes)\ z\ xs) \oplus (foldr\ (\otimes)\ z\ ys)\ .$$

We perform an induction on $xs$. The case when $xs := [\,]$ trivially holds. For $xs := x : xs$, we reason:

$$foldr\ (\otimes)\ (foldr\ (\otimes)\ z\ ys)\ (x : xs)$$
$$= x \otimes foldr\ (\otimes)\ (foldr\ (\otimes)\ z\ ys)\ xs$$
$$= \quad \{\,\text{induction}\,\}$$
$$x \otimes ((foldr\ (\otimes)\ z\ xs) \oplus (foldr\ (\otimes)\ z\ ys))$$
$$= \quad \{\,\text{assumption: } x \otimes (y \oplus w) = (x \otimes y) \oplus w\,\}$$
$$(x \otimes (foldr\ (\otimes)\ z\ xs)) \oplus (foldr\ (\otimes)\ z\ ys)$$
$$= (foldr\ (\otimes)\ z\ (x : xs)) \oplus (foldr\ (\otimes)\ z\ ys)\ .$$

Direction ($\Rightarrow$). Given $foldr\ (\otimes)\ z = hom\ (\oplus)\ (\otimes z)\ z$, prove that $x \otimes (y \oplus w) = (x \otimes y) \oplus w$ for $y$ and $w$ in the range of $foldr\ (\otimes)\ z$.

Let $y = foldr\ (\otimes)\ z\ xs$ and $w = foldr\ (\otimes)\ z\ ys$ for some $xs$ and $ys$. We reason:

$$x \otimes (y \oplus w)$$
$$= x \otimes (foldr\ (\otimes)\ z\ xs \oplus foldr\ (\otimes)\ z\ ys)$$
$$= \quad \{\,\text{since } foldr\ (\otimes)\ z = hom\ (\oplus)\ (\otimes z)\ z\,\}$$
$$x \otimes (foldr\ (\otimes)\ z\ (xs\!+\!\!+ys))$$
$$= foldr\ (\otimes)\ z\ (x : xs\!+\!\!+ys)$$
$$= \quad \{\,\text{since } foldr\ (\otimes)\ z = hom\ (\oplus)\ (\otimes z)\ z\,\}$$
$$foldr\ (\otimes)\ z\ (x : xs) \oplus foldr\ (\otimes)\ z\ ys$$
$$= (x \otimes foldr\ (\otimes)\ z\ xs) \oplus foldr\ (\otimes)\ z\ ys$$
$$= (x \otimes y) \oplus w\ .$$

$\square$