# Forward- or Reverse-Mode Automatic Differentiation: What's the Difference?

Birthe van den Berg[a,*], Tom Schrijvers[a], James McKinna[b], Alexander Vandenbroucke[c]

[a]*KU Leuven, Celestijnenlaan 200A, Leuven, Belgium*
[b]*Heriot-Watt University, Earl Mountbatten Building G.52, Edinburgh EH14 4AS, United Kingdom*
[c]*Standard Chartered, 1 Basinghall Avenue, London, EC2V 5DD, United Kingdom*

**Abstract**

Automatic differentiation (AD) has been a topic of interest for researchers in many disciplines, with increased popularity since its application to machine learning and neural networks. Although many researchers appreciate and know how to apply AD, it remains a challenge to truly understand the underlying processes. From an algebraic point of view, however, AD appears surprisingly natural: it originates from the differentiation laws. In this work we use Algebra of Programming techniques to reason about different AD variants, leveraging Haskell to illustrate our observations. Our findings stem from *three fundamental algebraic abstractions*: (1) the notion of semimodule, (2) Nagata's construction of the 'idealization of a module', and (3) Kronecker's delta function, that together allow us to write a *single-line abstract definition* of AD. From this single-line definition, and by instantiating our algebraic structures in various ways, we derive different AD variants, that have the same *extensional* behaviour, but different *intensional* properties, mainly in terms of (asymptotic) computational complexity. We show the different variants equivalent by means of *Kronecker isomorphisms*, a further elaboration of our Haskell infrastructure which guarantees correctness by construction. With this framework in place, this paper seeks to make AD variants more comprehensible, taking an algebraic perspective on the matter.

*Keywords:* symbolic differentiation, automatic differentiation, monoids, Cayley representation, semirings, Nagata idealization

*Corresponding author

*Email addresses:* `birthe.vandenberg@kuleuven.be` (Birthe van den Berg),
`tom.schrijvers@kuleuven.be` (Tom Schrijvers), `j.mckinna@hw.ac.uk` (James McKinna),
`alexander.vandenbroucke@sc.com` (Alexander Vandenbroucke)

## 1. Introduction

Many algorithms in artificial intelligence crucially depend not only on the efficient large-scale evaluation of arithmetic expressions, but also on the evaluation of their partial derivatives at specific points [1, 2]. Symbolic and numerical differentiation can be used for these computations but suffer from performance and accuracy issues, respectively. As a consequence, there is much interest in a third approach, called automatic differentiation (AD) [3]. AD efficiently and accurately evaluates the derivatives of any computation composed of basic arithmetic and exponential, logarithmic, and trigonometric functions.

The orthodox accounts of automatic differentiation, since its first introduction by Wengert [3], emphasise it as a whole-program transformation. It is coupled with a shift from evaluation of the primitive operations of the underlying language of numerical computation to those on the so-called *dual numbers* (due originally to Clifford [4], considerably predating their subsequent application to AD). The various AD modes (forward, reverse, mixed) then arise as further optimisations of this basic picture in terms of the efficient organisation and traversal of the underlying datastructures, typically exploiting mutability.

The vast literature on AD, which is impossible to survey here, covers extensions of Wengert's original approach in all possible directions: algorithmic improvements, language designs and extended expressivity, novel applications and much more. The contribution of this paper is not to extend this body of work even further, but rather to revisit the original algorithms from an algebraic perspective. Following Bird and de Moor's "Algebra of Programming" [5], we show how algebra can be used to (re-)derive existing AD algorithms from a basic executable specification, expressed in the purely functional programming language Haskell. The algebra involved consists of largely standard constructions (with one notable exception, namely Nagata's construction described in Section 3.2 below), which—because they yield isomorphic mathematical objects—have no *extensional* consequences. Yet, they do establish the functional equivalence between the different algorithmic variants. At the algorithmic level the various (isomorphic) changes of representation are precisely what allows us to exploit *intensional* aspects, specifically to change the computational complexity.

Specifically, this paper makes the following contributions.

- First, we show how symbolic differentiation and evaluation of symbolic expressions can both be seen algebraically as unique semiring homomorphisms from the free semiring. In terms of code, they are both instances of the same fold-like recursion scheme. This algebraic view enables us first to fuse the two homomorphisms to obtain basic forward mode AD, and then to recover symbolic differentiation as an instance of that (Section 2).

- Next, we abstract the AD algorithm in terms of three fundamental algebraic abstractions, whose specification and (equational) properties we realise via Haskell type classes: *modules* over a semiring, a generalization of dual numbers which we dub *Nagata numbers* (whose construction appears originally due to Nagata [6]), and of Kronecker's *delta* function. This

2

abstract algorithm readily instantiates to a more efficient form of forward mode AD that computes the whole gradient vector of partial derivatives in one go (Section 3).

- Then, we show how other instances of the abstract algorithm can be obtained by replacing the representation of the gradient using an isomorphism of such structure, using Elliott's method of type class morphisms [7]. This is demonstrated on a sparse vector representation (Section 4).

- Using the above methodology, we ultimately derive purely functional and imperative variants of reverse-mode AD: we successively improve the time complexity of scalar multiplication and gradient vector addition with representations based on homomorphisms, and switch to their equivalent mutable representation (Section 5).

- As an epilogue, we demonstrate that our approach readily admits common, useful extensions: additional primitive functions, direct computation with Nagata numbers to eliminate the overhead in building and evaluating symbolic expressions, higher derivatives, and **let**-sharing (Section 6).

- We provide an implementation[1] of the approach and different AD versions.

Finally, the paper discusses related work (Section 7) and concludes (Section 8).

We believe that this paper may be of interest to algebraists, (Haskell) programmers and AD experts alike. While the paper does its best to be accessible to these three groups, we acknowledge that this requires working through unfamiliar concepts or unfamiliar presentations of familiar concepts before the three strands come together. Appendix C supports readers who are not familiar with Haskell, briefly introducing type classes and library functions that are used throughout this work.

## 2. Setting the Stage

This section provides a basic executable definition of automatic differentiation (AD) based on a minimal setup. For readers who want a gentle introduction to the methodology used here, we refer to "Domain-Specific Languages of Mathematics" [8].

### 2.1. Symbolic Expressions and their Evaluation

In order to focus on the essential algorithmic mechanics, throughout most of this paper we pare down to a minimal programming language in which functions can be expressed: symbolic expressions *Expr v* that capture polynomials over variables $v$, that vary independently[2].

---

[1] https://github.com/birthevdb/Abstract-AD
[2] In Section 6.5 we extend the language of expressions with a **let**-construct to permit variables that may be bound to complex subexpressions.

| identity: | commutativity: |
|---|---|
| $zero \oplus x = x = x \oplus zero$ | $x \oplus y = y \oplus x$ |
| $one \otimes x = x = x \otimes one$ | $x \otimes y = y \otimes x$ |
| **associativity:** | **distributivity:** |
| $x \oplus (y \oplus z) = (x \oplus y) \oplus z$ | $x \otimes (y \oplus z) = (x \otimes y) \oplus (x \otimes z)$ |
| $x \otimes (y \otimes z) = (x \otimes y) \otimes z$ | $(x \oplus y) \otimes z = (x \otimes z) \oplus (y \otimes z)$ |
| **annihilator:** | $zero \otimes x = zero = x \otimes zero$ |

Figure 1: The commutative semiring laws.

```
data Expr v = Var v | Zero | One
            | Plus (Expr v) (Expr v) | Times (Expr v) (Expr v)
```

For example, for univariate expressions we can define a type $X$ with one variable. The symbolic representation of $x \times (x + 1)$ is captured in $example_1$.

```
data X = X
example₁ :: Expr X
example₁ = Times (Var X) (Plus (Var X) One)
```

We assign an overloaded meaning to symbolic expressions by mapping them to a *(commutative) semiring*. The 'semi' in semiring refers to the fact that a semiring is a ring that does not need to have an additive inverse. We omit it here because AD does not require it. Nevertheless, Section 6.1 shows that, if such an inverse exists, it can be accommodated. The commutativity requirement we impose is not essential, but simplifies the presentation. Section 6.2 extends the approach to non-commutative semirings, which is necessary to support, for instance, matrices. The overloading gives us the flexibility to interpret expressions in terms of both their plain value and of different AD flavors.

**Definition 1 (Semiring).** *A* semiring *is a set, together with two binary operators $\oplus$ ("addition") and $\otimes$ ("multiplication") and two distinguished elements zero and one. These should obey the laws of Figure 1: $\oplus$ and $\otimes$ should both be assocative and commutative and have neutral elements zero and one respectively. Moreover, $\otimes$ should distribute over $\oplus$ and have zero as annihilator.*

In Haskell, we capture the algebraic concept of a semiring in the *Semiring* type class, with instances for numeric types (e.g., *Int*, *Float*[3]) that respect the laws.

---

[3] If we ignore the limitations of floating point arithmetic.

```
class Semiring d where            instance Num a ⇒ Semiring a where
  zero :: d                         zero = 0
  one  :: d                         one  = 1
  (⊕) :: d → d → d                  (⊕) = (+)
  (⊗) :: d → d → d                  (⊗) = (∗)
```

There is one instance of semirings that is of special interest to us: that for symbolic expressions $Expr\ v$, quotiented by the laws.[4]

```
instance Semiring (Expr v) where
  zero = Zero
  one  = One
  (⊕)  = Plus
  (⊗)  = Times
```

What makes this instance special is that it is the *free* semiring instance, which gives rise to a fold-style *eval* function. This evaluator proceeds by structural recursion over the symbolic expressions to map them to their interpretation in a semiring $d$, given a mapping $var :: v → d$ for variables.

```
eval :: Semiring d ⇒ (v → d) → Expr v → d
eval var (Var x)      = var x
eval var Zero         = zero
eval var One          = one
eval var (Plus   e₁ e₂) = eval var e₁ ⊕ eval var e₂
eval var (Times  e₁ e₂) = eval var e₁ ⊗ eval var e₂
```

This way we can evaluate $example_1$ for $X = 5$ with the *Int* semiring.

```
> eval (λX → 5) example₁
30
```

*2.2. Homomorphisms*

The notion of free semiring, the associated *eval* function and its properties are special cases of generic results from the Algebra of Programming [5] which hold for any algebraic structure. Another related notion, which will turn out to be central to this paper, is that of *homomorphism*. Informally, this is a function that preserves algebraic structure. More formally,

**Definition 2 (Semiring homomorphism).** *A (semiring) homomorphism between two semirings $D_1$ and $D_2$ is a function $h :: D_1 → D_2$ that preserves the semiring structure:*

---

[4]Quotienting by the laws means that we should consider two (even structurally different) expressions equal if the laws say so (e.g., *Plus Zero (Var X) == Var X*).

$$h \; zero_{D_1} = zero_{D_2} \qquad\qquad\qquad h \; (x \oplus_{D_1} y) = h \; x \oplus_{D_2} h \; y$$
$$h \; one_{D_1} \;\; = one_{D_2} \qquad\qquad\qquad h \; (x \otimes_{D_1} y) = h \; x \otimes_{D_2} h \; y$$

Semiring homomorphisms compose: given homomorphisms $h_1 :: D_1 \rightarrow D_2$ and $h_2 :: D_2 \rightarrow D_3$, their composition $h_2 \circ h_1 :: D_1 \rightarrow D_3$ is also a homomorphism; and the expected laws hold. Furthermore, the identity function $id :: D \rightarrow D$ is trivially a homomorphism. The function $eval \; var$ is a homomorphism for any choice of $var$. In fact, it is the *unique* homomorphism $h$ such that $h \; (Var \; v) = var \; v$. From this uniqueness property follow two useful properties:

**Property 1 (Fusion).** *The* fusion property *allows us to incorporate any homomorphism $h$ into evaluation:* $h \circ eval \; var = eval \; (h \circ var)$.

**Property 2 (Reflection).** *The* reflection property *observes that evaluation with Var is the identity:* $eval \; Var = id$.

*2.3. Symbolic Differentiation and Dual Numbers*

The *eval* homomorphism also allows us to do symbolic differentiation. We capture its rules in the *derive* function, using the *Semiring* instance of *Expr v*, and exploiting the commutativity in the multiplication case.

$$
\begin{array}{ll}
derive :: Eq \; v \Rightarrow v \rightarrow Expr \; v \rightarrow Expr \; v \\
derive \; x \; (Var \; y) & = \textbf{if } x == y \textbf{ then } one \textbf{ else } zero \\
derive \; x \; Zero & = zero \\
derive \; x \; One & = zero \\
derive \; x \; (Plus \;\;\; e_1 \; e_2) & = derive \; x \; e_1 \oplus derive \; x \; e_2 \\
derive \; x \; (Times \; e_1 \; e_2) & = (\boxed{e_2} \otimes derive \; x \; e_1) \oplus (\boxed{e_1} \otimes derive \; x \; e_2)
\end{array}
$$

These do not fit *eval*'s fold-style recursion scheme because, in the case for *Times*, $e_1$ and $e_2$ appear outside of recursive calls ( $\boxed{gray}$ ). Using the so-called "banana-split" property [9, 10, 11], we can turn this function into the proper shape: we make it return a tuple $(Expr \; v, Expr \; v)$ of the original expression and the symbolic derivative.

$$derive' \; x \; e = (e, derive \; x \; e)$$

We thus obtain the following structurally recursive definition for *derive'*, from which we may then recover *derive* by projection on the second component:

$$
\begin{array}{ll}
derive' :: Eq \; v \Rightarrow v \rightarrow Expr \; v \rightarrow (Expr \; v, Expr \; v) \\
derive' \; x \; (Var \; y) & = (Var \; y, \textbf{if } x == y \textbf{ then } one \textbf{ else } zero) \\
derive' \; x \; Zero & = (zero, \;\;\; zero) \\
derive' \; x \; One & = (one, \;\;\; zero) \\
derive' \; x \; (Plus \;\;\; e_1 \; e_2) & = \textbf{let } (e_1', de_1) = derive' \; x \; e_1 \\
& \qquad\quad\;\; (e_2', de_2) = derive' \; x \; e_2 \\
& \quad \textbf{in } \; (e_1' \oplus e_2', de_1 \oplus de_2) \\
derive' \; x \; (Times \; e_1 \; e_2) & = \textbf{let } (e_1', de_1) = derive' \; x \; e_1
\end{array}
$$

$$(e_2', de_2) = derive' \; x \; e_2$$
$$\textbf{in} \;\; (e_1 \otimes e_2', (e_2' \otimes de_1) \oplus (e_1' \otimes de_2))$$

Before we write the functionality of *derive'* in terms of *eval*, we generalize $(Expr \; v, Expr \; v)$ to tuples whose components are drawn from an arbitrary semiring. In the AD literature, these are known as *dual numbers*, originally due to Clifford [4].

$$\textbf{data} \; Dual \; d = D \; \{ \, pri^D :: d,$$
$$tan^D :: d \, \}$$

**instance** *Functor Dual* **where**
$$fmap :: (a \to b) \to Dual \; a \to Dual \; b$$
$$fmap \; h \; (D \; f \; df) = D \; (h \; f) \; (h \; df)$$

Dual numbers are functorial; *fmap f* transforms both parameters using $f$. The two components of a dual number $D \; x \; dx$ are respectively known as the *primal* $x$ and the *tangent dx*. For a semiring $d$, the type *Dual d* also admits a semiring structure, where the definition of tangent precisely follows the corresponding cases of *derive*.

**instance** *Semiring* $d \Rightarrow$ *Semiring* (*Dual d*) **where**
$$zero \qquad\qquad\quad = D \; zero \quad zero$$
$$one \qquad\qquad\quad = D \; one \quad\; zero$$
$$(D \; f \; df) \oplus (D \; g \; dg) = D \; (f \oplus g) \; (df \oplus dg)$$
$$(D \; f \; df) \otimes (D \; g \; dg) = D \; (f \otimes g) \; ((g \otimes df) \oplus (f \otimes dg))$$

Now, symbolic differentiation may be seen as another instance of *eval*uation, but into the semiring of dual numbers over symbolic expressions.

$$symbolic :: Eq \; v \Rightarrow v \to Expr \; v \to Dual \; (Expr \; v)$$
$$symbolic \; x = eval \; gen \; \textbf{where} \; gen \; y = D \; (Var \; y) \; (\delta_x \; y)$$
$$\delta_x \quad y = \textbf{if} \; x == y \; \textbf{then} \; one \; \textbf{else} \; zero$$

*2.4. Classic Forward Automatic Differentiation*

Automatic differentiation (AD) is the umbrella term for algorithms that programmatically compute both the value *and* the derivative of an expression *at a point* [12, 13, 14]. Thus, if we combine symbolic differentiation with subsequent evaluation into a semiring $d$, we obtain a specification for forward-mode AD:

$$forwardAD \; var \; x = fmap \; (eval \; var) \circ symbolic \; x$$

We can make this specification more efficient by exploiting *eval*'s fusion property.

$$fmap \; (eval \; var) \circ symbolic \; x$$
$$= \;\; \{\text{- definition of } symbolic \text{ -}\}$$
$$fmap \; (eval \; var) \circ eval \; gen \; \textbf{where} \; gen \; y = D \; (Var \; y) \; (\delta_x \; y)$$
$$\delta_x \; y \;\; = \textbf{if} \; x == y \; \textbf{then} \; one \; \textbf{else} \; zero$$
$$= \;\; \{\text{- } eval \text{ fusion (Property 1) -}\}$$
$$eval \; (fmap \; (eval \; var) \circ gen) \; \textbf{where} \; gen \; y = D \; (Var \; y) \; (\delta_x \; y)$$
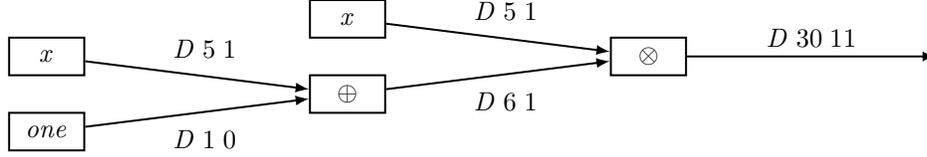$$\delta_x \; y \;\; = \textbf{if} \; x == y \; \textbf{then} \; one \; \textbf{else} \; zero$$

Figure 2: Running forward-mode AD on $example_1$ at $X = 5$.

After inlining $fmap$ ($eval\ var$) into $gen$, we get the classic definition of forward-mode AD, which calls $eval$ only once. Differences compared with $symbolic$ are highlighted in gray.

$$forwardAD :: (Eq\ v, Semiring\ d) \Rightarrow (v \to d) \to v \to Expr\ v \to Dual\ d$$
$$forwardAD\ \boxed{var}\ x = eval\ gen\ \textbf{where}\ gen\ y = D\ (\boxed{var}\ y)\ (\delta_x\ y)$$
$$\delta_x\ y\ = \textbf{if}\ x == y\ \textbf{then}\ one\ \textbf{else}\ zero$$

For example, the derivative of $example_1$ at the point $X = 5$ is 11; the process computing the result is depicted in Figure 2.

$$> forwardAD\ (\lambda X \to 5)\ X\ example_1$$
$$D\ 30\ 11$$

*2.5. Symbolic AD as Forward AD*

Observe that the only essential difference between symbolic and automatic differentiation is that the former evaluates into the symbolic $Expr\ v$ semiring, while the latter uses a 'numeric' semiring $d$. In fact, we can also have $forwardAD$ use the symbolic semiring, interpreting variables with $Var$.

$$forwardAD\ Var$$
$$=\ \ \{\text{- AD specification (Specification 2.4) -}\}$$
$$fmap\ (eval\ Var) \circ symbolic$$
$$=\ \ \{\text{- reflection property (Property 2) -}\}$$
$$fmap\ id \circ symbolic$$
$$=\ \ \{\text{- functor law and } id \text{ as unit of composition -}\}$$
$$symbolic$$

In summary, $symbolic$ is a special case of $forwardAD$:

$$symbolic = forwardAD\ Var$$

This sets the stage for the remainder of this paper: we use algebraic techniques to reason about the structure of differentiation algorithms and by exposing more abstract structure, one and the same abstract algorithm (Section 3) subsumes various concrete AD flavors.

$$
\begin{array}{lll}
e \oplus_{\mathrm{E}} zero_{\mathrm{E}} & = e & \qquad zero_{\mathrm{D}} \quad \bullet\, e = zero_{\mathrm{E}} \\
e \oplus_{\mathrm{E}} (e' \oplus_{\mathrm{E}} e'') = (e \oplus_{\mathrm{E}} e') \oplus_{\mathrm{E}} e'' & \qquad (d \oplus_{\mathrm{D}} d') \bullet e = (d \bullet e) \oplus_{\mathrm{E}} (d' \bullet e) \\
e \oplus_{\mathrm{E}} e' & = e' \oplus_{\mathrm{E}} e & \qquad one_{\mathrm{D}} \quad \bullet\, e = e \\
& & \qquad (d \otimes_{\mathrm{D}} d') \bullet e = d \bullet (d' \bullet e) \\[1ex]
d \bullet zero_{\mathrm{E}} \quad = zero_{\mathrm{E}} \\
d \bullet (e \oplus_{\mathrm{E}} e') = (d \bullet e) \oplus_{\mathrm{E}} (d \bullet e')
\end{array}
$$

Figure 3: The laws for a $D$-module $E$ over a semiring $D$, where $E$ is a commutative monoid. We overload the monoidal operation $\oplus$ for the corresponding *Semiring*, to emphasize its additive interpretation.

## 3. Abstract Automatic Differentiation

The central claim of this paper is that we can express the two previous differentiation functions (symbolic differentiation and classic forward-mode AD) as well as several more variants to come later, as the same one-line generic program. This section identifies three fundamental algebraic structures to accomplish this: $d$-modules $e$, generalized dual numbers $d \ltimes e$, and the Kronecker delta function. It also shows how to instantiate these to recover a version of forward-mode AD that computes the gradient vector of all partial derivatives *at once*, rather than for a single variable at a time.

*3.1. Modules over Semirings*

Our first abstraction generalizes the notion of gradient $\nabla e$ of an expression $e$.

$$
\nabla e = \left[ \begin{array}{c} \frac{\partial e}{\partial x_1} \\ \vdots \\ \frac{\partial e}{\partial x_n} \end{array} \right]
$$

By analogy with vector spaces (over a field) or modules (over a ring), we consider semimodules over a semiring $d$ [15], which we refer to as *d-modules*.

**Definition 3 (Semimodule).** *Given a semiring $D$, then a $D$-module $E$ is defined as a commutative monoid $E$, together with* scalar multiplication ($\bullet$) :: $D \to E \to E$, *which distributes over the additive structure of $E$. This means that $E$ should have a commutative, associative binary operator $\oplus$ with neutral element zero. Figure 3 gives $E$'s monoid laws and the laws that govern the interaction between the monoid structure of $E$ and the scalar multiplication.*

In Haskell, we capture the concept of a semimodule in the type class *Module*,[5] on top of the *Monoid* type class.

---

[5] We use the functional dependency $e \to d$ to anchor a $d$-module $e$ to its underlying semiring $d$ (Appendix C).

```
class Monoid e where                    class (Semiring d, Monoid e)
    zero :: e                                ⇒ Module d e | e → d where
    (⊕) :: e → e → e                         (•) :: d → e → e
```

Every semiring $d$ may be viewed trivially as a $d$-module, by taking $(•) = (⊗)$. A well-known example that is not itself a semiring is the vector space $\mathbb{R}^3$; it forms a monoid with vector addition and the zero vector $\begin{bmatrix} 0 & 0 & 0 \end{bmatrix}$, and an $\mathbb{R}$-module with the usual scalar multiplication.

*3.2. Nagata Numbers*

The second abstraction generalizes dual numbers of Section 2.3 to two components $d$ and $e$ that have *distinct* types, reflecting the idea that the value of an expression (its primal), and that of its derivative (tangent), conceptually live in different spaces. We exploit precisely the generalization that permits primal and tangent *to vary independently*. This concept first appeared in the work of Nagata [6] in the early 1950s, under the name "idealization of a module". However, the application of Nagata's construction to the study of AD appears to be novel. By analogy with dual numbers, we call the type $d ⋉ e$ *Nagata numbers*.

```
data d ⋉ e = N { priᴺ :: d,              instance Functor (d ⋉ ·) where
                 tanᴺ :: e }                 fmap h (N f df) = N f (h df)
```

Observe that the type constructor $d ⋉ ·$ is functorial in its second parameter and a *bi*functor in both parameters (Appendix A).

Nagata proved the following fundamental theorem.

---

**Theorem 1.** *Given a $d$-module $e$, then $d ⋉ e$ admits a semiring structure.*

---

The Haskell type class instance below witnesses this theorem.

```
instance Module d e ⇒ Semiring (d ⋉ e) where
    zero                   = N zero    zero
    one                    = N one     zero
    (N f df) ⊕ (N g dg) = N (f ⊕ g) (df ⊕ dg)
    (N f df) ⊗ (N g dg) = N (f ⊗ g) ((f • dg) ⊕ (g • df))
```

*3.3. Kronecker Delta*

The final piece of our jigsaw comes from the observation that the function $δ_x\ y = \textbf{if } x == y \textbf{ then } one \textbf{ else } zero$ is the Kronecker delta function.

$$\delta_x(y) = \begin{cases} 1 & x = y \\ 0 & otherwise \end{cases}$$

This gives rise to our third fundamental algebraic abstraction. In Haskell, we capture a generalized notion of this concept in the *Kronecker* type class.

```
class Module d e ⇒ Kronecker v d e where
    delta :: v → e
```

Notice that, at this point, no extra laws are imposed on *Kronecker v d e*; later, we will use isomorphism conditions (Definition 5) to constrain different AD variants. The idea is that *delta x* does not compute $\delta_x(y)$ for just a single variable $y$, but for *all* variables at once. Thus *delta x* is a value in a $d$-module $e$ where only the $x$-component is *one*, and all others are *zero*. A helpful analogy may be to think of $v$ as the type of basis vectors of $e$. Then, *delta x* is the embedding of basis vector $x$ into $e$. Continuing Section 3.1's example of the $\mathbb{R}$-module $\mathbb{R}^3$, there are usually three variables $x, y, z$ such that

$$
\begin{aligned}
delta\ x &= \begin{bmatrix} 1 & 0 & 0 \end{bmatrix} \\
delta\ y &= \begin{bmatrix} 0 & 1 & 0 \end{bmatrix} \\
delta\ z &= \begin{bmatrix} 0 & 0 & 1 \end{bmatrix}
\end{aligned}
$$

In the AD literature, these basis vectors, where one entry is *one* and all the other are *zero*, are often referred to as *one-hot vectors*.

### 3.4. Abstract Automatic Differentiation

With the three algebraic structures in place, we present our abstract definition of differentiation.

```
abstractD :: Kronecker v d e ⇒ (v → d) → Expr v → d ⋉ e
abstractD var = eval gen where gen x = N (var x) (delta x)
```

Next, we show how this abstract definition can be instantiated to recover *forwardAD*. Later, we vary $e$ to obtain more efficient versions.

### 3.5. Base case: Forward-Mode AD as AD in the Dense Function Space

We recover *forwardAD*—the classical forward-mode AD implementation given in Section 2.4—as an instance of *abstractD*, by using the standard function space $e = v → d$. We call this space *dense*[6] to contrast with the *sparse* function space coming up in Section 4.3, by analogy with dense/sparse polynomial representations in computer algebra [16].

```
type Dense v d = v → d
```

Informally, we can think of this structure as a functional representation of a gradient vector, required by least square fitting, gradient descent and many other applications. Just like a gradient vector, it contains one $d$-value per variable of type $v$. The actions of its monoid and $d$-module instances are given pointwise, using the '$d$-as-$d$-module construction' of Section 3.1.

---

[6]To avoid clutter, we present *Dense* and later types as **type** synonyms in the paper, but our implementation uses **newtype**s for unambiguous type class resolution.

**instance** $Semiring\ d \Rightarrow Monoid\ (Dense\ v\ d)$ **where**
$\quad zero \quad = \lambda v \rightarrow zero$
$\quad f_1 \oplus f_2 = \lambda v \rightarrow f_1\ v \oplus f_2\ v$

**instance** $Semiring\ d \Rightarrow Module\ d\ (Dense\ v\ d)$ **where**
$\quad d \bullet f \quad = \lambda v \rightarrow d \bullet (f\ v) = \lambda v \rightarrow d \otimes (f\ v)$

The Kronecker delta function for dense functions is the textbook one.

**instance** $(Eq\ v, Semiring\ d) \Rightarrow Kronecker\ v\ d\ (Dense\ v\ d)$ **where**
$\quad delta\ v = \lambda w \rightarrow$ **if** $v == w$ **then** $one$ **else** $zero$

Using this $Kronecker$ instance and generalizing dual numbers into $d \ltimes (Dense\ v\ d)$, we obtain the following definition for forward-mode AD.

$forwardAD_{Dense} :: (Eq\ v, Semiring\ d) \Rightarrow (v \rightarrow d) \rightarrow Expr\ v \rightarrow d \ltimes (Dense\ v\ d)$
$forwardAD_{Dense} = abstractD$

This is $forwardAD$ with its parameters re-arranged and $Dual\ d$ generalized:

$pri^D\ (forwardAD\ var\ x\ e) = pri^N\ (forwardAD_{Dense}\ var\ e\ x)$
$tan^D\ (forwardAD\ var\ x\ e) = tan^N\ (forwardAD_{Dense}\ var\ e\ x)$

The essential difference between $forwardAD$ and $forwardAD_{Dense}$ is intensional. In moving from computing a value in $v \rightarrow (d \ltimes d)$ to one in $d \ltimes (v \rightarrow d)$, we no longer compute both the primal and tangent for each individual $v$, but instead only compute the primal *once* for the whole gradient. In other words, we share the computation of the primal among all partial derivatives. For instance, $example_2$ represents expression $(x \times y) + x + 1$ in variables $x$ and $y$ using our expression language.

**data** $XY = X \mid Y$
$example_2 :: Expr\ XY$
$example_2 = Plus\ (Plus\ (Times\ (Var\ X)\ (Var\ Y))\ (Var\ X))\ One$

We compute its gradient at the point $(5, 3)$:

$> $**let** $\{var\ X = 5; var\ Y = 3; e = tan^N\ (forwardAD_{Dense}\ var\ example_2)\}$
$\quad$ **in** $(e\ X, e\ Y)$
$(4, 5)$

In contrast, when using $forwardAD$, we would have to call it twice:

$> $**let** $\{var\ X = 5; var\ Y = 3, e = example_2\}$
$\quad$ **in** $(tan^D\ (forwardAD\ var\ X\ e), tan^D\ (forwardAD\ var\ Y\ e))$

## 4. AD Variants by Construction

This section explains our methodology, which is based on constructing *Kronecker isomorphisms*, and demonstrates it on the sparse variant of $forwardAD_{Dense}$. This methodology will allow us to instantiate our abstract algorithm by varying the tangent type to obtain different flavors of automatic differentiation.

### 4.1. Relating AD Variants with Kronecker Isomorphisms

A key concern when devising new AD variants is their correctness. We can establish this by showing their functional equivalence to an established AD baseline. We do so using *Kronecker* isomorphisms.

**Definition 4 (Kronecker Homomorphism).** *Given a type of variables $V$ and a semiring $D$, with types $E_1$ and $E_2$ such that Kronecker $V$ $D$ $E_1$ and Kronecker $V$ $D$ $E_2$ hold (and thus also Monoid $E_1$, Monoid $E_2$, Module $D$ $E_1$ and Module $D$ $E_2$), then a* Kronecker homomorphism *is a function $h :: E_1 \rightarrow E_2$ that preserves the Kronecker $V$ $D$ structure:*

$$
\begin{aligned}
h \; zero_{E_1} &= zero_{E_2} & h \; (d \bullet_{E_1} m) &= d \bullet_{E_2} (h \; m) \\
h \; (x \oplus_{E_1} y) &= (h \; x) \oplus_{E_2} (h \; y) & h \; (delta_{E_1} \; v) &= delta_{E_2} \; v
\end{aligned}
$$

**Definition 5 (Kronecker Isomorphism).** *A* Kronecker isomorphism *is a pair of Kronecker homomorphisms $h :: E_1 \rightarrow E_2$ and $h^{-1} :: E_2 \rightarrow E_1$, that are inverses: $h \circ h^{-1} = id$ and $h^{-1} \circ h = id$.*

In the rest of the paper we use *Dense $V$ $D$* as the baseline. Hence, when coming up with a new tangent type $E$, we also want to find a Kronecker isomorphism $h :: Dense \; V \; D \rightarrow E$ to establish its correctness. This isomorphism is a function pair consisting of a representation function *rep* and its inverse *abs* (abstraction), following the terminology of data representation [17, 18]. To streamline our approach, we capture this isomorphism in a type class *CorrectAD*.

```
class Kronecker v d e ⇒ CorrectAD v d e where
    rep :: Dense v d → e
    abs :: e → Dense v d
```

With the *free* theorem [19] that follows from the parametricity of *abstractD* [20] we can then relate the new instance of the algorithm to the baseline.

---

**Theorem 2.** *Given a CorrectAD $V$ $D$ $E$ instance, we have that*

$$
\begin{aligned}
fmap \; abs \; (abstractD_E \; var \; e) &= abstractD_{Dense} \; var \; e \\
abstractD_E \; var \; e &= fmap \; rep \; (abstractD_{Dense} \; var \; e)
\end{aligned}
$$

---

### 4.2. Kronecker Isomorphisms by Construction

Given the $Kronecker\ V\ D\ (Dense\ V\ D)$ structure as our starting point and a newly desired representation $E$, we do not have to devise the $Kronecker\ V\ D\ E$ structure and the Kronecker isomorphism out of the blue. Instead, following the method outlined by Elliott [7], we can constructively derive both from a plain invertible function $h :: Dense\ V\ D \to E$. Indeed, we can create the instances for $E$ in terms of the instances of $Dense\ V\ D$ and of $h$.

**instance** $Monoid\ E$ **where**
$\quad zero\ = h\ zero_{\text{Dense}}$
$\quad x \oplus y = h\ (h^{-1}\ x \oplus_{\text{Dense}} h^{-1}\ y)$
**instance** $Module\ D\ E$ **where**
$\quad d \bullet x = h\ (d \bullet_{\text{Dense}} h^{-1}\ x)$

**instance** $Kronecker\ V\ D\ E$ **where**
$\quad delta\ v = h\ (delta_{\text{Dense}}\ v)$
**instance** $CorrectAD\ V\ D\ E$ **where**
$\quad rep\ = h$
$\quad abs\ = h^{-1}$

By construction, these instances both (1) satisfy all the necessary laws, and (2) ensure that $h$ is a Kronecker isomorphism. For example, we can verify that $x \oplus zero = x$ and that $h\ x \oplus h\ y = h\ (x \oplus y)$ as follows:

$\quad x \oplus zero$
$= \quad \{$- definitions of $zero$ and $\oplus$ -$\}$
$\quad h\ (h^{-1}\ x \oplus h^{-1}\ (h\ zero))$
$= \quad \{$- isomorphism -$\}$
$\quad h\ (h^{-1}\ x \oplus zero)$
$= \quad \{$- monoid identity law for $Dense\ v\ d$ -$\}$
$\quad h\ (h^{-1}\ x)$
$= \quad \{$- isomorphism -$\}$
$\quad x$

$\quad h\ x \oplus h\ y$
$= \quad \{$- definition of $\oplus$ -$\}$
$\quad h\ (h^{-1}\ (h\ x) \oplus h^{-1}\ (h\ y))$
$= \quad \{$- isomorphism -$\}$
$\quad h\ (x \oplus y)$

Verifying the other laws and homomorphism properties proceeds similarly.

In short, each time we devise a new AD variant, we come up with a new representation $E$ and an invertible function $h :: Dense\ V\ D \to E$. The rest of the infrastructure then follows mechanically, except for the fact that we can simplify the definitions, sometimes exploiting further equations, for the sake of optimizing the (algorithmic properties of the) representation.

The invertible function ensures, through the constructed Kronecker isomorphism, that the extensional behaviour remains the same, and thus guarantees correctness. At the same time, a well-chosen alternative representation may improve intensional properties like asymptotic runtime.

### 4.3. Sparse Maps

We demonstrate the above approach to switch from the dense function space representation to that of sparse maps. This exploits the fact that many partial derivatives of sub-expressions are zero. Indeed, if a variable $x$ does not occur in an expression $e$, the partial derivative $\frac{\partial e}{\partial x}$ is zero. Notably, in the base instance of dense functions, when the expression is just a variable, all partial derivatives but

one are by definition zero. We avoid explicit representations of such zeros—as well as explicit operations with them—by replacing the dense function space $v \rightarrow d$ with the type of finite key-value maps $Map\ v\ d$. We use Haskell's $Data.Map$ library (Appendix C) instead of a specialized dictionary representation [21].

**type** $Sparse\ v\ d = Map\ v\ d$

The omitted values in the map are $zero$. This way we recover the $Dense$ representation from the sparse one by looking up the variables in the map and defaulting to $zero$ when a variable is not present.

$$abs_{\rightarrow} \quad :: (Ord\ v, Semiring\ d) \Rightarrow Sparse\ v\ d \rightarrow Dense\ v\ d$$
$$abs_{\rightarrow}\ m \quad = \lambda v \rightarrow findWithDefault\ zero\ v\ m$$

$$rep_{\rightarrow} \quad :: (Ord\ v, Bounded\ v, Enum\ v, Semiring\ d, Eq\ d)$$
$$\Rightarrow Dense\ v\ d \rightarrow Sparse\ v\ d$$
$$rep_{\rightarrow}\ dense = fromList\ [(v, dense\ v) \mid v \leftarrow [minBound \mathbin{.\,.} maxBound],$$
$$dense\ v \not\equiv zero]$$

Observe that $abs_{\rightarrow}$ and $rep_{\rightarrow}$ are inverse functions under the constraints (modulo the equivalence of explicit and implicit zeroes). With the new $Sparse\ v\ d$ representation and the two inverse functions we can derive the $Monoid$, $Module$, $Kronecker$ and $CorrectAD$ instances following the approach outlined in Section 4.1. This results in the all-zero map being empty. Addition of maps is performed key-wise, but avoided if either or both maps have an implicit zero at a given (variable) key.

**instance** $(Ord\ v, Semiring\ d) \Rightarrow Monoid\ (Sparse\ v\ d)$ **where**
$\quad zero = empty$
$\quad (\oplus)\ = unionWith\ (\oplus)$

Scalar multiplication is also key-wise. As $d \otimes zero = zero$ (Figure 1), we can effortlessly preserve the implicit zero entries.

**instance** $(Ord\ v, Semiring\ d) \Rightarrow Module\ d\ (Sparse\ v\ d)$ **where**
$\quad d \bullet m = fmap\ (d\otimes)\ m$

The only non-zero entry for $delta\ x$ is that for $x$.

**instance** $(Ord\ v, Semiring\ d) \Rightarrow Kronecker\ v\ d\ (Sparse\ v\ d)$ **where**
$\quad delta\ x = singleton\ x\ one$

Moreover, and crucially in what follows, we observe the following relationship between this $Kronecker$ structure and its $Module$ and $Monoid$ instances:

$$d \bullet delta\ v = singleton\ v\ d \tag{1}$$
$$e \oplus (d \bullet delta\ v) = insertWith\ (\oplus)\ v\ d\ e \tag{2}$$

Lastly, the $CorrectAD$ instance is straightforward.

$$\textbf{instance}\ (Ord\ v, Bounded\ v, Enum\ v, Semiring\ d, Eq\ d) \Rightarrow$$
$$CorrectAD\ v\ d\ (Sparse\ v\ d)\ \textbf{where}$$
$$rep = rep_\rightarrow$$
$$abs = abs_\rightarrow$$

In summary, all that the sparse gradient computation requires, in comparison to the dense computation, is the switch from $d \ltimes (Dense\ v\ d)$ to $d \ltimes (Sparse\ v\ d)$.

$$forwardAD_{Sparse} :: (Ord\ v, Semiring\ d) \Rightarrow (v \to d) \to Expr\ v \to d \ltimes \boxed{Sparse\ v\ d}$$
$$forwardAD_{Sparse} = abstractD$$

Theorem 2 now holds by construction. We can verify this correspondence by revisiting $example_2$ at the point $(5, 3)$.

$$> \textbf{let}\ \{var\ X = 5; var\ Y = 3\}\ \textbf{in}\ tan^N\ (forwardAD_{Sparse}\ var\ example_2)$$
$$\{X \mapsto 4; Y \mapsto 5\}$$

## 5. Reverse-Mode Automatic Differentiation

We now develop the appropriate $d$-module structure that exhibits reverse-mode AD as an instance of $abstractD$ by refining the above sparse representation. Motivated by the quest for improved efficiency, we optimize the time complexity of the essential functions for abstract differentiation with the underlying Nagata numbers semiring: scalar multiplication ($\bullet$), addition ($\oplus$) and Kronecker's $delta$.

### 5.1. Accumulating Multiplications

We start with another function-based representation. Instead of the $d$-module $e$, we consider the space of *$d$-module homomorphisms* from $d$ to $e$. It is a standard result in commutative algebra [22, p.18][7] that these two structures are isomorphic, considered as $d$-modules. In terms of the *underlying* Haskell types, this amounts to taking the plain function space.

$$\textbf{type}\ d \multimap e = d \to e$$

But we intend this type $d \multimap e$ to represent only the $d$-module homomorphisms. This is the case for *linear* functions from $d$ to $e$. A function $f :: d \to e$ is linear[8] iff it has the multiplicative *homogeneity* property:

$$f\ (x \otimes y)\ =\ x \bullet f\ y$$

As the Haskell type $d \to e$ also includes non-linear functions, we have to make disciplined use of $d \multimap e$, being careful to avoid the non-linear ones. Intuitively,

---

[7] For the case of commutative *rings* and their modules; the argument for semirings and modules is identical, but simpler, because there are fewer equations to check.

[8] Preservation of addition follows from the fact that the result of $f$ is a semimodule.

we can think of $d \multimap e$ as augmenting $e$ with an accumulator for a scalar multiplier of type $d$. By supplying the scalar $one$, we recover the original module.

$$rep_{\multimap} :: Module\ d\ e \Rightarrow e \to (d \multimap e) \qquad abs_{\multimap} :: Module\ d\ e \Rightarrow (d \multimap e) \to e$$
$$rep_{\multimap}\ e = \lambda d \to d \bullet e \qquad\qquad abs_{\multimap}\ f = f\ one$$

In particular, as $one \bullet e = e$ (Figure 3), we have that $abs_{\multimap} \circ rep_{\multimap} = id$. Moreover, for linear functions also $rep_{\multimap} \circ abs_{\multimap} = id$. Now we derive $Monoid$, $Module$ and $Kronecker$ instances that make $rep_{\multimap}$ a homomorphism by construction. It follows that $d \multimap e$ inherits its monoidal structure in a pointwise fashion from $e$. For $e = Dense\ v\ d$, the time complexity for additive operations remains linear in the number of variables.

> **instance** $Monoid\ e \Rightarrow Monoid\ (d \multimap e)$ **where**
> $\quad zero\ \ = \lambda d \to zero$
> $\quad f \oplus g = \lambda d \to f\ d \oplus g\ d$

Furthermore, its $d$-module structure uses the function parameter as an accumulator for the scalar multiplier. This affords a *"strength reduction"*: the expensive scalar multiplication ($\bullet$) is replaced by the cheaper semiring multiplication ($\otimes$).

> **instance** $Module\ d\ e \Rightarrow Module\ d\ (d \multimap e)$ **where**
> $\quad d' \bullet f = \lambda d \to f\ (d' \otimes d)$

This reduces the time complexity of scalar multiplication from $\mathcal{O}(V)$ (where $V$ denotes the number of variables with non-zero values in the given sparse map) for the $Sparse\ v\ d$ representation to $\mathcal{O}(1)$ for that of $d \multimap Sparse\ v\ d$. However, the general $Kronecker$ instance of $d \multimap e$ still contains the expensive ($\bullet$) operator.

> **instance** $Kronecker\ v\ d\ e \Rightarrow Kronecker\ v\ d\ (d \multimap e)$ **where**
> $\quad delta\ v = \lambda d \to d \bullet delta\ v$

Fortunately, we only require a specialized instance for sparse maps, which operates in constant time, as a consequence of equation (1) for $d \bullet delta\ v$.

> **instance** $(Ord\ v, Semiring\ d) \Rightarrow Kronecker\ v\ d\ (d \multimap Sparse\ v\ d)$ **where**
> $\quad delta = singleton$

Finally, the $CorrectAD$ instance for $d \multimap e$ is defined as the composition of $rep_{\multimap}/abs_{\multimap}$ with $rep/abs$ of $e$.

> **instance** $CorrectAD\ v\ d\ e \Rightarrow CorrectAD\ v\ d\ (d \multimap e)$ **where**
> $\quad rep = rep_{\multimap} \circ rep$
> $\quad abs = abs \circ abs_{\multimap}$

In summary, by using $d \ltimes (d \multimap Sparse\ v\ d)$ we obtain reverse-mode AD.

$$reverseAD :: (Ord\ v, Semiring\ d) \Rightarrow (v \to d) \to Expr\ v \to d \ltimes (\boxed{d \multimap Sparse\ v\ d})$$
$$reverseAD = abstractD$$

For instance, consider the more extensive example expression $x \times ((x+1) \times (x+x))$, encoded in Haskell as $example_3$.

$$example_3 :: Expr\ X$$
$$example_3 = Times\ (Var\ X)\ (Times\ (Plus\ (Var\ X)\ One)\ (Plus\ (Var\ X)\ (Var\ X)))$$

Running reverse-mode AD on this expression, the derivative at $X = 5$ is 170.

$$> (abs_{\multimap} \circ tan^N)\ (reverseAD\ (\lambda X \to 5)\ example_3)$$
$$\{X \mapsto 170\}$$

Figure 4 shows that this computation effectively works in multiple passes, which is a key characteristic of reverse-mode AD. During the forward (bottom-up) pass (shown by the black arrows) the algorithm computes Nagata numbers, which contain the value of each sub-expression, and the dependencies between the nodes as function closures $c :: d \multimap Sparse\ v\ d$, e.g., $N\ 300\ c_{\otimes_2}$ for the final outcome. In the backward (top-down) pass (shown by the red left-to-right arrows), we invoke the function closure $c_{\otimes_2}$ on argument $one$ (using $abs_{\multimap}$). It propagates the scalar accumulator downwards. At each downward step, the accumulator changes to reflect how much the overall derivative changes given that the subexpression changes by one. Finally, the results at the leaves are added (shown by the blue right-to-left arrows) to yield the overall derivative. The conventional imperative algorithms typically accomplish this with imperative updates. We derive that form in Section 5.3, but first we address the linear cost of addition.

*5.2. Accumulating Additions*

We can optimize the additive structure of $d$-modules in much the same way as we have optimized scalar multiplications: by considering a representation of $e \to e$ functions that have the *additive* homogeneity property:

$$f\ (x \oplus y) = x \oplus f\ y$$

This representation is also known as the Cayley representation [23, 24] of $e$, the most well-known example of which is difference lists [25] or Hughes lists [26]. Again, the Haskell type $e \to e$ includes functions that do not have the homogeneity property; those are not included in the *Cayley e* type.

**type** $Cayley\ e = e \to e$

$rep_{\diamond} :: Monoid\ e \Rightarrow e \to Cayley\ e$      $abs_{\diamond} :: Monoid\ e \Rightarrow Cayley\ e \to e$
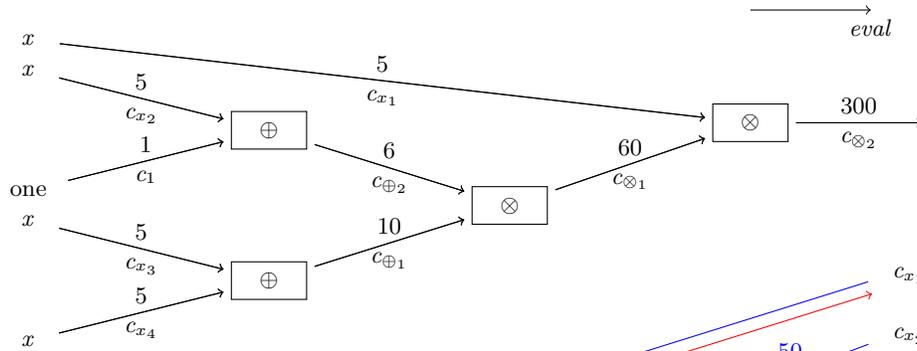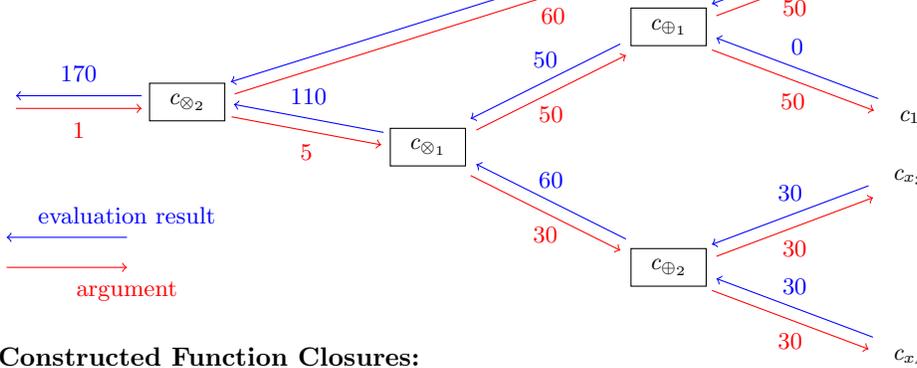$rep_{\diamond}\ e = \lambda e' \to e' \oplus e$      $abs_{\diamond}\ f = f\ zero$

This representation typically implies another "strength reduction", namely from $e$'s binary operator ($\oplus$) to function composition ($\circ$). We thereby reduce the time complexity of addition from $\mathcal{O}(V)$ for $Sparse\ v\ d$ to $\mathcal{O}(1)$ for $Cayley\ (Sparse\ v\ d)$.

**Forward Pass:**



**Backward Pass:**



**Constructed Function Closures:**

| Derivative | Closure | Argument ($d$) | Result |
|---|---|---|---|
| $c_{\otimes_2}$ | $\lambda d \to c_{x_1}\ (d \otimes 60) \oplus c_{\otimes_1}\ (d \otimes 5)$ | 1 | 170 |
| $c_{\otimes_1}$ | $\lambda d \to c_{\oplus_1}\ (d \otimes 10) \oplus c_{\oplus_2}\ (d \otimes 6)$ | 5 | 110 |
| $c_{\oplus_1}$ | $\lambda d \to c_{x_2}\ d \oplus c_1\ d$ | 50 | 50 |
| $c_{\oplus_2}$ | $\lambda d \to c_{x_3}\ d \oplus c_{x_4}\ d$ | 30 | 60 |
| $c_1$ | $\lambda d \to \{\,\}$ | 50 | 0 |
| $c_{x_i}$ | $\lambda d \to \{\, x \mapsto d\,\}$ | 60, 50, 30, 30 | 60, 50, 30, 30 |

Figure 4: Reverse-mode AD of $x \times (x + 1) \times (x + x)$, for $x = 5$.

19

> **instance** *Monoid* (*Cayley e*) **where**
>    *zero*   = *id*
>    $f \oplus g = g \circ f$

We derive the *Module* and *Kronecker* instances following our standard recipe. However, we cannot eliminate the expensive operations from the *Cayley e* instances for *Module* and *Kronecker*. Indeed, the *Module* instance requires converting from *Cayley e* to *e* and back, where notably $rep_\diamond$ re-introduces *e*'s expensive $\mathcal{O}(V)$ addition operation ($\oplus$), which we have just managed to eliminate from the *Monoid* instance.

> **instance** *Module d e* $\Rightarrow$ *Module d* (*Cayley e*) **where**
>    $d \bullet f = rep_\diamond \, (d \bullet abs_\diamond \, f)$

The *Kronecker* instance derived via $rep_\diamond$ from *e* similarly features *e*'s addition operation ($\oplus$).

> **instance** *Kronecker v d e* $\Rightarrow$ *Kronecker v d* (*Cayley e*) **where**
>    *delta v* $= \lambda e \to e \oplus delta\ v$

Fortunately, we can again define a specialized *Kronecker* instance for the type $e = d \multimap Cayley$ (*Sparse v d*), which avoids appeal to the expensive $\oplus$ operation. Indeed, the ($\bullet$) operation defined on *e* need not make use of that defined on *Cayley* (*Sparse v d*). For the specialized *delta*, we instead calculate:

> *delta v* =   {- definition of *delta* for $d \multimap$ -}
>       $\lambda d \to d \bullet delta\ v$
>   =   {- definition of ($\bullet$) and *delta* for *Cayley* -}
>       $\lambda d \to rep_\diamond \, (d \bullet abs_\diamond \, (\lambda e \to e \oplus delta\ v))$
>   =   {- definition of $rep_\diamond$ and $abs_\diamond$ -}
>       $\lambda d \to \lambda e \to e \oplus (d \bullet (\lambda e \to e \oplus delta\ v)\ zero)$
>   =   {- $\beta$-reduction and unit of monoid (Figure 3) -}
>       $\lambda d \to \lambda e \to e \oplus (d \bullet delta\ v)$
>   =   {- equation (2) -}
>       $\lambda d \to \lambda e \to insertWith\ (\oplus)\ v\ d\ e$
>   =   {- $\eta$-reduction -}
>       $\lambda d \to insertWith\ (\oplus)\ v\ d$

Hence, we obtain an implementation of *delta* that takes only $\mathcal{O}(\log V)$ time.

> **instance** (*Ord v, Semiring d*) $\Rightarrow$ *Kronecker v d* ($d \multimap Cayley$ (*Sparse v d*))
>   **where** *delta* = *insertWith* ($\oplus$)

Moreover, the *Cayley e* representation preserves the correctness of *e*.

> **instance** *CorrectAD v d e* $\Rightarrow$ *CorrectAD v d* (*Cayley e*) **where**
>    $rep = rep_\diamond \circ rep$
>    $abs = abs \circ abs_\diamond$

In summary, with type $d \multimap Cayley\ (Sparse\ v\ d)$, both ($\bullet$) and ($\oplus$) take $\mathcal{O}(1)$ and $delta$ takes $\mathcal{O}(\log V)$. As a consequence, the time complexity of $abstractD$ for a program with $N$ nodes and $V$ variables is now $\mathcal{O}(N \times \log V)$.

$$reverseAD_{Cayley} :: (Ord\ v, Semiring\ d)$$
$$\Rightarrow (v \to d) \to Expr\ v \to d \ltimes (\ \boxed{d \multimap Cayley\ (Sparse\ v\ d)}\ )$$
$$reverseAD_{Cayley} = abstractD$$

For example, as before, the derivative of $example_3$ at $X = 5$ is 170.

$$> (abs_\diamond \circ abs_{\multimap} \circ tan^N)\ (reverseAD_{Cayley}\ (\lambda X \to 5)\ example_3)$$
$$\{X \mapsto 170\}$$

Notice that we can see values of type $Cayley\ (Sparse\ v\ d)$ as a functional model for computations that modify a state of type $Sparse\ v\ d$, i.e., immutable maps.

### 5.3. Array-based Reverse Mode

Performance-wise, we can do even better by switching from immutable maps to mutable arrays (from $Data.Array.ST$) (Appendix C). Using these to represent the state, instead of the purely functional implementation above, we can perform in-place updates. Furthermore, the array can be provided by an implicit environment ($Control.Monad.Reader$). Hence, instead of type $Cayley\ (Sparse\ v\ d)$, we can work with the type $ReaderT\ (STArray\ s\ v\ d)\ (ST\ s)\ ()$. It is convenient to encapsulate this rather wordy type in the type synonym $STCayley\ v\ d$.

$$\textbf{type}\ STCayley\ v\ d = \forall s.ReaderT\ (STArray\ s\ v\ d)\ (ST\ s)\ ()$$

Observe that $STCayley\ v\ d$ polymorphically quantifies over the state thread parameter $s$, which is necessary if we want to run the computation using Haskell's built-in primitive $runST :: (\forall s.ST\ s\ a) \to a$.

The isomorphism between $Cayley\ (Sparse\ v\ d)$ and $STCayley\ v\ d$ essentially converts between the array and map representations.

$$rep_\gg :: (Ix\ v, Semiring\ d) \Rightarrow Cayley\ (Sparse\ v\ d) \to STCayley\ v\ d$$
$$rep_\gg\ p = foldMap\ (\lambda(v, d) \to modifyAt\ (\oplus d)\ v)\ (toList\ (p\ zero))$$
$$abs_\gg :: (Ix\ v, Bounded\ v, Semiring\ d) \Rightarrow STCayley\ v\ d \to Cayley\ (Sparse\ v\ d)$$
$$abs_\gg\ q = runST\ (\textbf{do}\ arr \leftarrow newArray\ (minBound, maxBound)\ zero$$
$$runReaderT\ q\ arr$$
$$l \leftarrow getAssocs\ arr$$
$$return\ \$\ foldMap\ (\lambda(v, d) \to insertWith\ (\oplus)\ v\ d)\ l)$$

The $rep_\gg$ function uses a helper function $modifyAt$ to edit the value at a specific position of the array, which is the counterpart of $insertWith$ on maps.

$$modifyAt :: Ix\ v \Rightarrow (d \to d) \to v \to STCayley\ v\ d$$
$$modifyAt\ f\ v = \textbf{do}\ arr \leftarrow ask; x \leftarrow readArray\ arr\ v; writeArray\ arr\ v\ (f\ x)$$

Using the methodology of Section 4.1 the *Monoid* instance works out to be the generic instance that is available to any type $m$ () where $m$ is a monad.

> **instance** *Monoid* (*STCayley v d*) **where**
>   *zero*  = *return* ()
>   $p \oplus q = p \gg q$

As we work with $d \multimap STCayley\ v\ d$, the *Module* structure of *STCayley v d* is not used. For that reason, we skip immediately to the *Kronecker* instance where the purpose of using mutable arrays is to reduce the $\mathcal{O}(\log V)$ time complexity of *insertWith* in *delta* to the $\mathcal{O}(1)$ time complexity of *modifyAt*.

> **instance** (*Ix v, Semiring d*) $\Rightarrow$ *Kronecker v d* (*STCayley v d*) **where**
>   *delta v = modifyAt* ($\oplus one$) *v*

In the $d \multimap STCayley\ v\ d$ combination this becomes:

> **instance** (*Ix v, Semiring d*) $\Rightarrow$ *Kronecker v d* ($d \multimap STCayley\ v\ d$) **where**
>   *delta v* $= \lambda d \rightarrow$ *modifyAt* ($\oplus d$) *v*

The correctness of *STCayley v d* is established in terms of that of *Cayley* (*Sparse v d*).

> **instance** (*Ix v, Ord v, Bounded v, Enum v, Semiring d, Eq d*) $\Rightarrow$
>   *CorrectAD v d* (*STCayley v d*) **where**
>     $rep = rep_{\gg} \circ rep$
>     $abs = abs \circ abs_{\gg}$

Finally, we get the resulting monadic definition for reverse-mode AD.

> $reverseAD_m :: (Ix\ v, Semiring\ d)$
>               $\Rightarrow (v \rightarrow d) \rightarrow Expr\ v \rightarrow d \ltimes ( \boxed{d \multimap STCayley\ v\ d} )$
> $reverseAD_m = abstractD$

We use this efficient version of reverse-mode AD on $example_3$ as follows:

> $> (abs_{\diamond} \circ abs_{\gg} \circ abs_{\multimap} \circ tan^N)\ (reverseAD_m\ (\lambda X \rightarrow 5)\ example_3)$
> $\{X \mapsto 170\}$

### 5.4. Summary

Figure 5 summarizes the successive differentiation algorithms we have constructed by instantiating *abstractD* with semiring types $d \ltimes e$ for different choices of tangent type $e$. It also shows the worst-case time complexity of each algorithm. Notice that, when $V$ is large, i.e., there are many variables, reverse-mode AD, and the last two variants in particular, is the most interesting. However, for few variables, forward-mode AD is preferable as it only takes a single pass.

Finally, in Figure 6 we collapse the layers of abstraction to show what the one-but-last entry of Figure 5—our most efficient purely functional reverse-mode AD algorithm—looks like after aggressive inlining. For the sake of simplicity we have also replaced the custom algebraic datatype for Nagata numbers with plain Haskell tuples.

| Differentiation Mode | Semiring Type | Time Complexity |
|---|---|---|
| Symbolic | $(Expr\ v) \ltimes (Expr\ v)$ | $\mathcal{O}(N^2)$ |
| Forward Mode | $d \ltimes (Dense\ v\ d)$ | $\mathcal{O}(N \times V)$ |
| Forward Mode | $d \ltimes (Sparse\ v\ d)$ | $\mathcal{O}(N \times V)$ |
| Reverse Mode | $d \ltimes (d \multimap Sparse\ v\ d)$ | $\mathcal{O}(N \times V)$ |
| Reverse Mode | $d \ltimes (d \multimap Cayley\ (Sparse\ v\ d))$ | $\mathcal{O}(N \times \log V)$ |
| Reverse Mode | $d \ltimes (d \multimap STCayley\ v\ d)$ | $\mathcal{O}(N + V)$ |

Figure 5: The successive refinements of the semiring type, with their corresponding differentiation modes and time complexity, for an expression of $N$ nodes in $V$ variables.

$$
\begin{aligned}
&autodiff :: (Ord\ v, Semiring\ d) \Rightarrow (v \to d) \to Expr\ v \to Map\ v\ d\\
&autodiff\ var\ e = snd\ (go\ var\ e)\ one\ zero\ \textbf{where}\\
&\quad go\ :: (Ord\ v, Semiring\ d)\\
&\qquad \Rightarrow (v \to d) \to Expr\ v \to (d, d \to Map\ v\ d \to Map\ v\ d)\\
&\quad go\ var\ (Var\ v) \quad\ = (var\ v, \lambda d \to \lambda m \to insertWith\ (\oplus)\ v\ d\ m)\\
&\quad go\ var\ Zero \qquad\ = (zero,\ \ \lambda d \to \lambda m \to m)\\
&\quad go\ var\ One \qquad\ \ = (one,\ \ \lambda d \to \lambda m \to m)\\
&\quad go\ var\ (Plus\quad e_1\ e_2) = \textbf{let}\ (f, df) = go\ var\ e_1\\
&\qquad\qquad\qquad\qquad\quad (g, dg) = go\ var\ e_2\\
&\qquad\qquad\qquad\quad \textbf{in}\ \ (f \oplus g, \lambda d \to \lambda m \to df\ d\ (dg\ d\ m))\\
&\quad go\ var\ (Times\ e_1\ e_2) = \textbf{let}\ (f, df) = go\ var\ e_1\\
&\qquad\qquad\qquad\qquad\quad (g, dg) = go\ var\ e_2\\
&\qquad\qquad\qquad\quad \textbf{in}\ \ (f \otimes g, \lambda d \to \lambda m \to df\ (d \otimes g)\ (dg\ (d \otimes f)\ m))
\end{aligned}
$$

Figure 6: Fully inlined AD code based on the $d \ltimes (d \multimap Cayley\ (Sparse\ v\ d))$ semiring.

## 6. Extensions

Our setup allows for various extensions that make AD more useful. Notice that, in these extensions, we often revert to the more naive $forwardAD_{Dense}$ with the *Dense* representation for the sake of brevity.

### 6.1. Additional Primitives

We have so far focused only on the minimal *Semiring* interface. Yet, it is well-known that dual numbers (and, by extension, Nagata numbers) easily adapt to other primitives, such as trigonometric, exponential and logarithmic functions. We model such functions as subclasses of the *Semiring* type class. For instance, we can extend the approach from semirings to rings by creating a *Ring* subclass with a *negate* method for additive inverses.[9]

---

[9]This deviates from the Haskell Prelude where *negate* is a method in the *Num* type class.

$$\textbf{class } \textit{Semiring } d \Rightarrow \textit{Ring } d \qquad \textbf{instance } (\textit{Module } d\ e, \textit{Ring } d) \Rightarrow \textit{Ring } (d \ltimes e)$$
$$\textbf{where } \textit{negate} :: d \rightarrow d \qquad \textbf{where } \textit{negate } (N\ f\ df) = N\ (-f)\ (-\textit{one} \bullet df)$$

Here, we use the fact that Haskell allows us to write $-x$ for $\textit{negate } x$.

Likewise, we collect trigonometric functions in *Trig*, whose Nagata number instance relies on the well-known chain rule. In general, to enrich $d$ (and hence *Expr v*) with differentiable primitives requires us to know both the operation ($sin$), *and* its derivative ($cos$). Here, in order to be able to negate $sin$ as the derivative of $cos$, $d$ must be a *Ring*.

$$\textbf{class } \textit{Trig } d \textbf{ where} \qquad \textbf{instance } (\textit{Module } d\ e, \textit{Trig } d, \textit{Ring } d) \Rightarrow \textit{Trig } (d \ltimes e)$$
$$\textit{sin} :: d \rightarrow d \qquad \textbf{where } \textit{sin } (N\ f\ df) = N\ (\textit{sin } f)\ (\quad \textit{cos } f \bullet df)$$
$$\textit{cos} :: d \rightarrow d \qquad \textit{cos } (N\ f\ df) = N\ (\textit{cos } f)\ (-\textit{sin } f \bullet df)$$

Observe that, while $d$ must now also have a *Trig* instance, no additional demands are made of $e$. Hence, all the AD variations we have covered, which are only based on variations of $e$, also support additional primitive functions.

*6.2. Non-Commutative Semirings*

We can generalize our approach from commutative semirings to semirings with non-commutative multiplication ($\otimes$), such as (square) matrices over semirings and the semiring of formal languages over an alphabet.

The key change in our approach is that rather than starting from the following rule for symbolic differentation, which exploits the commutativity:

$$\textit{derive } x\ (\textit{Times } e_1\ e_2) = (\ \boxed{e_2 \otimes}\ \ \textit{derive } x\ e_1) \oplus (\ \boxed{e_1 \otimes}\ \ \textit{derive } x\ e_2)$$

where $e_1$ and $e_2$'s derivatives are *both scaled on the left*, we instead use the more conventional rule:

$$\textit{derive } x\ (\textit{Times } e_1\ e_2) = (\textit{derive } x\ e_1\ \boxed{\otimes e_2}\ ) \oplus (\ \boxed{e_1 \otimes}\ \ \textit{derive } x\ e_2)$$

which scales $\textit{derive } x\ e_2$ on the left and $\textit{derive } x\ e_1$ on the right.

More abstractly, we can model this with *d-d-bimodules* that feature distinct left-scaling and right-scaling operators.

$$\textbf{class } \textit{Semiring } d \Rightarrow \textit{Bimodule } d\ e \mid e \rightarrow d \textbf{ where}$$
$$(\overrightarrow{\bullet}) :: d \rightarrow e \rightarrow e$$
$$(\overleftarrow{\bullet}) :: e \rightarrow d \rightarrow e$$

Note that, for every $d$ whose multiplication is commutative, every *Module d e* instance gives rise to a *Bimodule d e* instance, by taking the two scaling operations to be the same underlying ($\bullet$) operation. If multiplication in $d$ is not commutative, nevertheless the distinguished case $e = d$ yields a *Bimodule* structure with ($\overrightarrow{\bullet}$) (resp. ($\overleftarrow{\bullet}$)) given by *left*-(resp. *right*-)multiplication (see below).

24

However, it is *not* sufficient to take (a restricted subset of) $d \to e$ as a representative of the set of *Bimodule* homomorphisms from $d$ to $e$.

Instead, we can avoid introducing this new $d$-$d$-bimodule abstraction and reuse the existing *Module*-related infrastructure by means of $(d, d)$-modules. At the core of this idea is the following semiring construction for pairs.

> **instance** *Semiring* $d \Rightarrow$ *Semiring* $(d, d)$ **where**
> $\quad zero = (zero, zero)$
> $\quad one \ = (one, one)$
> $\quad (lx, rx) \oplus (ly, ry) = (lx \oplus ly, rx \oplus ry)$
> $\quad (lx, rx) \otimes (ly, ry) = (lx \otimes ly, \boxed{ry \otimes rx}\ )$

This is the usual component-wise instance, with one exception: the factors in the multiplication are flipped for the second component. Then $d$ is a $(d, d)$ module as follows:

> **instance** *Semiring* $d \Rightarrow$ *Module* $(d, d)$ $d$ **where**
> $\quad (l, r) \bullet d = l \otimes d \otimes r$

Here, the first component of the tuple scales from the left and the second component scales from the right. We can make use of the left and right injections below to turn a $d$ value either into a left or right factor, using the neutral element *one* for the other factor.

> $inl, inr :: Semiring\ d \Rightarrow d \to (d, d)$
> $inl \ \ x = (x, one)$
> $inr \ \ y = (one, y)$

More generally, with a $(d, d)$-module $e$ we obtain a variation on the Nagata numbers:

> **instance** $(Semiring\ d, Module\ \boxed{(d, d)}\ e) \Rightarrow Semiring\ (Nagata\ d\ e)$ **where**
>
> $\quad zero = N\ zero\ zero$
> $\quad one \ = N\ one\ zero$
> $\quad N\ f\ df \oplus N\ g\ dg = N\ (f \oplus g)\ (df \oplus dg)$
> $\quad N\ f\ df \otimes N\ g\ dg = N\ (f \otimes g)\ ((\boxed{inr}\ g \bullet df) \oplus (\boxed{inl}\ f \bullet dg))$

This one localized modification seamlessly integrates non-commutative multiplication into our AD variants. In particular, following the earlier recipe we now use the $(d, d) \multimap e$ representation for accummulating multiplications, which is that of $(d, d) \to e$ functions that take two $d$ accumulators. The scalar multiplication of this representation updates both accumulators.

> **instance** $(Semiring\ d, Module\ (d, d)\ e) \Rightarrow Module\ (d, d)\ ((d, d) \multimap e)$ **where**
> $\quad (dl', dr') \bullet f = \lambda(dl, dr) \to f\ (dl' \otimes dl, dr \otimes dr')$

This shows that the first accumulator is for multiplications on the left and the other for multiplications on the right.

In short, $(d, d)$-modules support non-commutative multiplication with minimal upheaval.

### 6.3. Overloading Initial Expressions

Automatic differentiation avoids the intermediate symbolic derivative when computing precise numeric results. This is shown by the elimination of *eval* in Section 2.4, which was invoked twice in the naive implementation. In this subsection we show that, by overloading operators in the initial expression and thereby avoiding the initial symbolic expression, we can further eliminate the *eval* call and thus further reduce verbosity.

In particular, we assume that the initial expression is not given as a value of type *Expr v*, but as a polymorphic function of type $\forall d.\,Semiring\ d \Rightarrow d \rightarrow ... \rightarrow d \rightarrow d$ with as many inputs of type $d$ as there are variables $v$. For example, $example_2$ from Section 3.5, can be written as a polymorphic function of arity 2.

$$example'_2 :: Semiring\ d \Rightarrow d \rightarrow d \rightarrow d$$
$$example'_2\ x\ y = ((x \otimes y) \oplus x) \oplus one$$

The mechanism behind deriving $example'_2$ is that of *fold/build fusion* [27], applied to *eval var $example_2$*, where $example_2$ is an instance of the builder pattern, building a symbolic expression, and *eval* is a fold. Fold/build fusion eliminates the intermediate symbolic structure and directly builds the result of evaluation. Gibbons and Wu [28] show how to exploit such structure systematically (generalising to all algebraic structures defined by Haskell type classes, not only *Semirings*), by relating the 'deep' embedding in terms of symbolic expressions *Expr v* with the corresponding 'shallow' embedding $\forall d.\,Semiring\ d \Rightarrow (v \rightarrow d) \rightarrow d$. Here, *eval* is one half of an isomorphism pair between these two types, with the inverse given by instantiation at *Var*.

Hence, instead of the verbose calls (left), we readily invoke $example'_2$ (right).

```
> let { var X = 5; var Y = 3 }          > example'₂ 5 3
  in eval var example₂                  21
21
> let { var X = 5; var Y = 3 }          > example'₂ (N 5 (delta X))
  in forwardAD_Sparse var example₂                (N 3 (delta Y))
N 21 { X ↦ 4; Y ↦ 5 }                   N 21 { X ↦ 4; Y ↦ 5 }
```

Both cases avoid the interpretative overhead of *eval*uating the symbolic expression. Yet, we can also recover the original symbolic expression.

```
> example'₂ (Var X) (Var Y)
Plus (Plus (Times (Var X) (Var Y)) (Var X)) One
```

### 6.4. Higher Derivatives

Higher derivatives are obtained through repeated differentiation. For instance, the second-order derivative is the derivative of the derivative. The naive approach to compute the second-order derivative at a point is through repeated evaluation of the symbolic derivative.
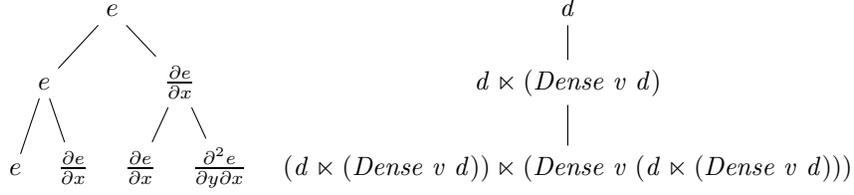
Figure 7: $0^{th}$–$2^{nd}$-order derivatives and their types.

$$derive_{2^{nd}} :: (Eq\ v, Semiring\ d) \Rightarrow (v \to d) \to Expr\ v \to Dense\ v\ (Dense\ v\ d)$$
$$derive_{2^{nd}}\ var\ e = \lambda x\ y \to eval\ var\ (tan^N\ (forwardAD_{Dense}\ Var$$
$$(tan^N\ (forwardAD_{Dense}\ Var\ e)\ x))\ y)$$

The Haskell call $derive_{2^{nd}}\ var\ e\ x\ y$ corresponds to the mathematical notation $\left.\frac{\partial}{\partial y}\left(\frac{\partial}{\partial x}e\right)\right|_{var}$. By using naturality of $tan^N$ and $eval$'s fusion property, we merge the three traversals of $derive_{2^{nd}}$ into a single $eval$ (see Appendix A):

$$derive'_{2^{nd}}\ var\ e = \lambda x\ y \to tan^N\ (tan^N\ (forwardAD_{Dense}\ gen\ e)\ x)\ y$$
$$\textbf{where}\ gen\ z = N\ (var\ z)\ (delta\ z)$$

This instance of $abstractD$ uses the $(d \ltimes (Dense\ v\ d))$-module of nested Nagata numbers of type $(d \ltimes (Dense\ v\ d)) \ltimes (Dense\ v\ (d \ltimes (Dense\ v\ d)))$.

The nested representation (Figure 7) of the second-order derivative contains redundancy: the value $\frac{\partial e}{\partial x}$ appears twice. This redundancy gets worse as the order (and thus the level of nesting) increases. In general, for the $n$th-order derivative we get $2^n$ values of which only $n + 1$ are distinct. To avoid this, it is more economical to use a compact representation.

$$\textbf{data}\ N_{2^{nd}}\ v\ d = N_{2^{nd}}\ d\ (Dense\ v\ (d \ltimes (Dense\ v\ d)))$$

The corresponding forward AD is similar to $abstractD$:

$$forwardAD_{2^{nd}} :: (Semiring\ d) \Rightarrow (v \to d) \to Expr\ v \to N_{2^{nd}}\ v\ d$$
$$forwardAD_{2^{nd}}\ var = eval\ gen\ \textbf{where}\ gen\ x = N_{2^{nd}}\ (var\ x)\ (delta\ x)$$

Running this on $example_3$ gives the evaluation result, first and second derivative:

$$> forwardAD_{2^{nd}}\ (\lambda X \to 5)\ example_3$$
$$N_{2^{nd}}\ 300\ (\lambda X \to N\ \{pri^N = 170, tan^N = (\lambda X \to 64)\})$$

This generalizes straightforwardly to other higher derivatives.

*All Higher Partial Derivatives.* In the extreme case where $n = \infty$, the economical representation becomes an infinite stream of all successive derivatives [29].

27

This stream of derivatives has the form $x :< dxs$ where $x$ is the value and $dxs$ are all higher-order partial derivatives with respect to all variables $v$.[10]

$$\textbf{data } Stream \ v \ d = d :< Dense \ v \ (Stream \ v \ d)$$

Then, the semiring operations for the infinite stream of successive derivatives are defined as follows.

$$\textbf{instance } Semiring \ d \Rightarrow Semiring \ (Stream \ v \ d) \ \textbf{where}$$
$$zero = zero :< zero$$
$$one \ = one :< zero$$
$$(x :< xs) \oplus (y :< ys) = (x \oplus y) :< (xs \oplus ys)$$
$$xxs \qquad \otimes yys \qquad = (x \otimes y) :< \lambda v \to ((xs \ v \otimes yys) \oplus (xxs \otimes ys \ v))$$
$$\textbf{where } xxs = (x :< xs)$$
$$yys = (y :< ys)$$

The *Module* and *Kronecker* instances follow naturally from this definition. They allow us to automatically differentiate a symbolic expression and obtain its stream of derivatives.

$$forwardAD_{All} :: (Eq \ v, Semiring \ d) \Rightarrow (v \to d) \to Expr \ v \to Stream \ v \ d$$
$$forwardAD_{All} \ var = eval \ gen \ \textbf{where } gen \ y = var \ y :< delta \ y$$
$$> forwardAD_{All} \ (\lambda X \to (5 :: Int)) \ example_3$$
$$300 :< (\lambda X \to 170 :< (\lambda X \to 64 :< (\lambda X \to 12 :< (\lambda X \to 0 :< ...$$

*Hybrid Automatic Differentiation.* In second-order optimization, it is often preferred to use a hybrid approach, i.e., combine forward mode and backward mode to get a more efficient algorithm. For instance, when taking a function's curvature into account, it is necessary to compute the Hessian[11] applied to a vector. Pearlmutter [31] introduces this mix of the two modes. In a forward sweep the first derivative is computed, whereupon the second, reverse sweep computes the second-order derivative at a point. This can be obtained by exploiting the $d$-module isomorphisms to work with

$$\textbf{data } N_H \ v \ d = N_H \ d \ (Dense \ v \ (d \ltimes (d \multimap (Cayley \ (Sparse \ v \ d)))))$$

### 6.5. Basic Sharing of Subexpressions

Our approach so far ignores any sharing in the source expression. Consider for instance $example_4$, which denotes the expression $(x+x) \times (x+x)$, but at the level of Haskell shares the two occurrences of $x + x$ by means of a **let** binding.

---

[10]It differs from Hinze's stream representation [30] in the interposition of *Dense* in the tail.

[11]The Hessian of a function $f$ in $n$ variables $x_i$ is the $n \times n$ square matrix consisting of all second-order derivatives $\frac{\partial^2 f}{\partial x_i \partial x_j}$.

$$example_4 :: Expr\ X$$
$$example_4 = \textbf{let}\ y = Plus\ (Var\ X)\ (Var\ X)\ \textbf{in}\ Times\ y\ y$$

The *eval* recursion scheme used by *abstractD* does not see the **let** binding and processes both occurrences of $x+x$. In particular, $abstractD\ (\lambda X \to 5)\ example_4$ reduces to $(N\ 5\ one \oplus N\ 5\ one) \otimes (N\ 5\ one \oplus N\ 5\ one)$ which features two $(\oplus)$ operations instead of one. To avoid this recomputation of intermediate results, we extend the symbolic expression language with a let-construct so that we can write $example'_4$ with this new syntax.

$$\textbf{data}\ Expr\ v = \ ...$$
$$\qquad\quad |\ Let\ v\ (Expr\ v)$$
$$\qquad\qquad\qquad (Expr\ v)$$

$$example'_4 :: Expr\ XY$$
$$example'_4 = Let\ Y\ (Plus\quad (Var\ X)\ (Var\ X))$$
$$\qquad\qquad\qquad\qquad (Times\ (Var\ Y)\ (Var\ Y))$$

We adapt the evaluator accordingly with the standard semantics for **let**-binding.

$$eval :: (Eq\ v, Semiring\ d) \Rightarrow (v \to d) \to Expr\ v \to d$$
$$...$$
$$eval\ gen\ (Let\ y\ e_1\ e_2) = \textbf{let}\ d_1 \quad = eval\ gen\ e_1$$
$$\qquad\qquad\qquad\qquad\qquad\ gen_2 = \lambda x \to \textbf{if}\ x == y\ \textbf{then}\ d_1\ \textbf{else}\ gen\ x$$
$$\qquad\qquad\qquad\qquad \textbf{in}\ \ eval\ gen_2\ e_2$$

Observe that, because $e_2$ has the additional **let**-bound variable in scope, we evaluate it with the extended generator $gen_2$, which maps the free variables to their given *gen* meaning and the new **let**-bound variable $y$ to the meaning $d_1$ of $e_1$. The definition of *abstractD* is as before, now with support for sharing.

*Sharing for Sparse Maps.* Now, $abstractD\ (\lambda X \to 5)\ example'_4$ reduces to an intermediate form with only one $(\oplus)$ operation, and if we commit for the tangent to the sparse module, this reduces further in a way that nicely preserves the sharing.

$$\textbf{let}\ y = N\ 5\ one \oplus N\ 5\ one\ \textbf{in}\ (y \otimes y)$$
$$\equiv\ \textbf{let}\ y = N\ 5\ \ \{X \mapsto 1\} \oplus N\ 5\ \{X \mapsto 1\}\ \textbf{in}\ (y \otimes y)$$
$$\equiv\ \textbf{let}\ y = N\ 10\ (fromList\ [(X, 2)])\qquad\quad\ \textbf{in}\ (y \otimes y)$$
$$\equiv\ N\ (10\ \times\ 10)\ (fmap\ (10\ \times)\ \{X \mapsto 2\} \oplus fmap\ (10\ \times)\ \{X \mapsto 2\})$$
$$\equiv\ N\ 100\ (\{X \mapsto 20\} \oplus \{X \mapsto 20\})$$
$$\equiv\ N\ 100\ \{X \mapsto 40\}$$

*Sharing for Functions.* However, if we use a function-based module representation, such as *Dense* and our three reverse-mode variants, then the $(\oplus)$ is duplicated again later in the evaluation. This happens because the $(\oplus)$ under the function binder is not reduced until the function is applied. When both references of the shared function are applied, the $(\oplus)$ in its body will be evaluated each time. The following shows this for dense functions because it is the simplest of the function-based representations.

$$\textbf{let } y = N\ 5\ one \oplus N\ 5\ one \textbf{ in } (y \otimes y)$$
$$\equiv \textbf{let } y = N\ 5\ \ (\lambda X \to 1) \oplus N\ 5\ (\lambda X \to 1) \qquad\qquad \textbf{in } (y \otimes y)$$
$$\equiv \textbf{let } y = N\ 10\ (\lambda X' \to (\lambda X \to 1)\ X' \oplus (\lambda X \to 1)\ X') \textbf{ in } (y \otimes y)$$
$$\equiv (N\ 10\ (\lambda X' \to (\lambda X \to 1)\ X' \oplus (\lambda X \to 1)\ X')) \otimes$$
$$\quad (N\ 10\ (\lambda X' \to (\lambda X \to 1)\ X' \oplus (\lambda X \to 1)\ X'))$$
$$\equiv N\ (10\ \times\ 10)\ (\lambda X''' \to (\lambda X'' \to 10\ \times\ (\lambda X' \to (\lambda X \to 1)\ X' \oplus$$
$$(\lambda X \to 1)\ X')\ X'')\ X''' +$$
$$(\lambda X'' \to 10\ \times\ (\lambda X' \to (\lambda X \to 1)\ X' \oplus$$
$$(\lambda X \to 1)\ X')\ X'')\ X''')$$

Fortunately, we can address the problem of sharing for function-representations with a custom symbolic rule for *Let* (Proof in Appendix B)[12].

---

**Lemma 1.**

$$\frac{\partial(\textbf{let } y\ =\ e_1 \textbf{ in } e_2)}{\partial x_i} = \left.\frac{\partial e_2}{\partial y}\right|_{y=e_1} \times \frac{\partial e_1}{\partial x_i} + \left.\frac{\partial e_2}{\partial x_i}\right|_{y=e_1}$$

---

For instance, the expression $example_4'$ has partial derivative:

$$\frac{\partial example_4'}{\partial x} = \left.\frac{\partial y \times y}{\partial y}\right|_{y=x+x} \times \frac{\partial x + x}{\partial x} + \left.\frac{\partial y \times y}{\partial x}\right|_{y=x+x}$$
$$= ((x+x)+(x+x)) \times (1+1) + 0$$

This formula (Lemma 1) is captured in a *Letin* type class with a method *letin*, which in turn is defined in terms of $abs :: e \to v \to d$ of the *CorrectAD* type class, that extracts the *v*-component of a *d*-module *e*.

```
class CorrectAD v d e ⇒ Letin v d e where
    letin :: v → e → e → e
    letin y de₁ de₂ = ((abs de₂ y) ● de₁) ⊕ de₂
```

We specialize the *eval*uator used by *abstractD* to use this formula, also specializing its type *d* to $d \ltimes e$.

```
eval :: (Eq v, Letin v d e) ⇒ (v → d ⋉ e ) → Expr v → d ⋉ e
...
eval gen (Let y e₁ e₂) =
    let N f₁ df₁  = eval gen e₁
        gen₂      = λx → if x == y then N f₁ (delta y) else gen x
        N f₂ df₂  = eval gen₂ e₂
    in  N f₂ (letin y df₁ df₂)
```

---

[12]We ignore the fact that some primitive functions may have non-differentiable points.

This evaluator uses the standard substitution behavior for the primal, but defers to the *letin* method of the *Letin* type class for the behavior for the tangent. If we inline and specialize *letin* for $d \multimap Cayley\ (Sparse\ v\ d)$, we get a tight definition that computes the maps produced by $de_1$ and $de_2$ only once.

> **instance** $(Ord\ v, Semiring\ d) \Rightarrow Letin\ v\ d\ (d \multimap Cayley\ (Sparse\ v\ d))$ **where**
> $\quad letin\ y\ de_1\ de_2 = \lambda n\ m \to$ **let** $m' = de_2\ n\ m$ **in** $de_1\ (m'\ !\ y)\ m'$

With the above definition, the **let**-sharing is properly preserved.

*Complications.* In order to get the main points across, we have ignored two complicating aspects above.

Firstly, the way the *Let* case of *eval* extends the generator *gen* results in a lookup time that is linear in the number ($L$) of nested **let**-bindings. A more efficient datastructure (a map or array) can be used instead to improve the lookup time complexity to $\mathcal{O}(\log L)$ or $\mathcal{O}(1)$.

Secondly, we have implicitly assumed that **let**-bound variables do not shadow other variables, be they free variables or other **let**-bound variables. To support shadowing, the partial derivative for the shadowed variable $y$—or at least, as much of it as has been accumulated so far in $m$—should be saved before running $dx_2$ and restored afterwards. Regardless, other than for debugging purposes, one may want to purge the partial derivative for the **let**-bound variable $y$ from $m'$ as it is no longer relevant. Both together are accomplished as follows.

> $letin\ y\ de_1\ de_2 = \lambda n\ m \to$ **let** $old = lookup\ y\ m$
> $\qquad\qquad\qquad\qquad\qquad\quad m' = de_2\ n\ m$
> $\qquad\qquad\qquad\qquad$ **in** $de_1\ (m'\ !\ y)\ (update\ (const\ old)\ y\ m')$

## 7. Related Work

There is a vast amount of literature on automatic differentiation since the idea was first presented by Wengert [3].

*State-of-the-Art AD Algorithms.* Much of the effort in recent years has been devoted to establishing the correctness of different AD variants, and to extending the supported language features. Most works focus on the operational semantics of formal calculi they define, and reason about the correctness of their AD approaches in operational terms. One of the most comprehensive and accessible of these, in our view, is the recent work of Krawiec et al. [12]. It too goes through several steps of refinement, from a standard forward mode AD to an efficient reverse mode variant. While we present several purely functional map-based variants, their main focus is the imperative version, which uses a defunctionalized form of our $d \multimap STCayley\ v\ d$ representation that resembles the so-called trace or Wengert list [3] used in many AD implementations. They cover several

additional language features, which we believe can also be incorporated into our approach, such as support for higher-order functions[13], sum and product types.

A virtue of their paper is in identifying criteria for a (successful!) approach to AD: that it should demonstrate efficient reverse-mode AD; that it be higher-order; that it be asymptotically efficient; and that it be provably correct. The approach we have presented here, thanks to our careful algebraic treatment, ensures that (functional) correctness is more or less immediate, via our use of Kronecker type-class isomorphisms; efficiency relative to the computation of the primal follows from our use of a single abstract computation *abstractD*, based on *eval*; and our ability to capture all the various modes of AD may be seen as arising by careful instantiation of Nagata's fundamental construction, which nevertheless may be summarised in completely elementary terms by the *Semiring* typeclass instance given in Section 3.2.

The authors of [12] give only informal arguments for the relation between different variants, but the correctness proof of their final version constitutes a non-trivial effort involving a Kripke logical relation. This proof approach is inspired by the earlier work of Huot et al. [32, 33] on a logical-relation-based correctness proof of a basic forward mode algorithm for a higher-order language.

Another recent work in a similar vein is that of Smeding and Vákár [34], who derive several variants of reverse-mode AD. They use "well-known program transformations", reasoning mostly in algorithmic terms.

*AD and Functional Programming.* Karczmarczuk [29] shows how to use type classes to overload arithmetic operations to obtain forward-mode AD with dual numbers, and generalizes that to an infinite stream of higher derivatives.

Elliott takes a more algebraic approach based on program derivation. First, he explores a derivation of forward-mode AD, focusing on the chain rule as the guiding principle, and generalizing from scalars to vector spaces [35]. Next, he shifts his view to derivatives as linear maps and the categorical structure of differentiation [36]. More recently, Elliott has also explored the application of automatic differentiation to languages [37].

Elsman et al. [38] follow Elliott's lead, with their functional analysis for differentiation based on linear functions. They show that symbolic and forward-mode derivatives follow from the same rules, but reverse-mode essentially computes the adjoint of a derivative. Adjoints arise by taking a curried form of $d$-valued inner product structure on vector spaces $e$, giving rise to a duality between $e$ and the linear function space $e \multimap d$. Our more general setting does not consider such normed/inner product structure on our $d$-modules $e$, but we expect that such richer structure may be smoothly incorporated into a general picture, and we leave this as a topic for future work.

Wang et al. [39] show that the two passes of reverse-mode AD can be conveniently implemented using delimited control in an impure functional setting

_____
[13]Observe that features like higher-order functions are parametric with respect to the semiring type and can thus be used as is.

like Scala; they provide no formal argument for its correctness.

Kmett's `ad` library [40] is a Haskell implementation based on Pearlmutter and Siskind's work [41], which explained how to incorporate AD as a first-class function in a functional language. The library exploits Haskell's abstraction and overloading mechanisms to support different AD modes.

Oleksandr et al. [42] discuss the challenges in automatic differentiation of higher-order functions, for example functions in curried form.

*Nagata Numbers.* As far as we know, our use of Nagata's idealization of a module ("Nagata numbers") in the context of AD is novel. Most works in the AD literature use the standard homogenous dual numbers. Some, like those of Krawiec et al. [12] and Smeding and Vákár [34], do feature various heterogeneous dual numbers, but they do not observe that these are instances of a general structure. Nagata originally defined his idealization of a module over a ring. Our weakening to consider only *semi*ring structure is just cosmetic, and the extension of our methods to consider a (commutative) *ring* $d$ of coefficients is straightforward, and we appeal to it without comment in Section 6. See also [43] for a comprehensive survey.

*Leibniz vs. Newton.* One way to understand our contribution, in contrast to those of Elliott and others, is that while they take a 'Newtonian' view of AD as an algebraic operation on spaces of (abstract) *functions*, we take the 'Leibnizian' approach of understanding both AD and *symbolic* differentiation as an algebraic operation on a language *Expr v* of (concrete) symbolic *expressions*. This in particular contrasts with Baydin et al.'s assertion [2, Section 2.2, under "What AD Is Not"] that AD is "not symbolic differentiation", but echoes their observation, following Griewank [44], that it has a "two-sided nature that is partly symbolic and partly numerical".

## 8. Conclusion

Having identified three algebraic abstractions, we can write symbolic differentiation, forward-mode and reverse-mode AD as different instances of one and the same abstract algorithm. A significant contribution is the observation that we can generalize the well-known dual numbers into the Nagata idealization abstraction $d \ltimes e$ for any $d$-module $e$, which allows the $e$ parameter to vary independently of the semiring $d$ of coefficients, and that we can use abstract Kronecker deltas to represent partial derivatives. Furthermore, guided by structure-preserving maps in the form of Kronecker isomorphisms, we have taken successive steps in optimising the representation, in a correct-by-construction manner, from basic symbolic differentiation via forward-mode AD to the (typically) much more efficient reverse-mode AD.

Future work includes further exploring matrix differentiation with this framework. Square matrices form a (non-commutative) semiring but matrices of arbitrary dimensions form an extra challenge, in particular when reflecting these dimensions in the types.

## 9. Acknowledgements

## References

[1] M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard, M. Kudlur, J. Levenberg, R. Monga, S. Moore, D. G. Murray, B. Steiner, P. A. Tucker, V. Vasudevan, P. Warden, M. Wicke, Y. Yu, X. Zhang, Tensorflow: A system for large-scale machine learning, CoRR abs/1605.08695 (2016). `arXiv:1605.08695`.

[2] A. G. Baydin, B. A. Pearlmutter, A. A. Radul, J. M. Siskind, Automatic differentiation in machine learning: A survey, J. Mach. Learn. Res. 18 (1) (2017) 5595–5637.

[3] R. E. Wengert, A simple automatic derivative evaluation program, Commun. ACM 7 (8) (1964) 463–464. `doi:10.1145/355586.364791`.

[4] W. K. Clifford, Preliminary sketch of bi-quaternions, in: Proceedings of the London Mathematical Society, Vol. 4, 1873, pp. 381—-395.

[5] R. Bird, O. de Moor, Algebra of Programming, Prentice-Hall, Inc., USA, 1997.

[6] M. Nagata, Local rings, Wiley Interscience, 1962.

[7] C. Elliott, Denotational design with type class morphisms (extended version), Tech. Rep. 2009-01, LambdaPix (March 2009).
URL `http://conal.net/papers/type-class-morphisms`

[8] P. Jansson, C. Ionescu, J.-P. Bernardy, Domain-Specific Languages of Mathematics, Vol. 24 of Texts in Computing, College Publications, 2022.

[9] G. Hutton, A tutorial on the universality and expressiveness of fold, J. Funct. Program. 9 (4) (1999) 355–372. `doi:10.1017/s0956796899003500`.

[10] E. Meijer, M. Fokkinga, R. Paterson, Functional programming with bananas, lenses, envelopes and barbed wire, in: Proceedings of the 5th ACM Conference on Functional Programming Languages and Computer Architecture, Springer-Verlag, Berlin, Heidelberg, 1991, p. 124–144.

[11] E. Meijer, Calculating compilers, Ph.D. thesis, Katholieke Universiteit Nijmegen (1992).

[12] F. Krawiec, S. Peyton Jones, N. Krishnaswami, T. Ellis, R. A. Eisenberg, A. W. Fitzgibbon, Provably correct, asymptotically efficient, higher-order reverse-mode automatic differentiation, Proc. ACM Program. Lang. 6 (POPL) (2022) 1–30. `doi:10.1145/3498710`.

[13] L. Szirmay-Kalos, Higher order automatic differentiation with dual numbers, Periodica Polytechnica Electrical Engineering and Computer Science 65 (1) (2021) 1–10. `doi:10.3311/PPee.16341`.

[14] P. H. Hoffmann, A hitchhiker's guide to automatic differentiation, Numerical Algorithms 72 (3) (2016) 775–811. `doi:10.1007/s11075-015-0067-6`.

[15] J. S. Golan, Semimodules over Semirings, Springer Netherlands, Dordrecht, 1999, pp. 149–161. `doi:10.1007/978-94-015-9333-5_14`.

[16] J. von zur Gathen, J. Gerhard, Modern Computer Algebra, 3rd edition, Cambridge University Press, 2013. `doi:10.1017/CBO9781139856065`.

[17] A. Gill, G. Hutton, The worker/wrapper transformation, Journal of Functional Programming 19 (2) (2009) 227–251. `doi:10.1017/S0956796809007175`.

[18] C. A. R. Hoare, Proof of correctness of data representations, Acta Informatica 1 (4) (1972) 271–281. `doi:10.1007/BF00289507`.

[19] P. Wadler, Theorems for free!, in: J. E. Stoy (Ed.), Proceedings of the fourth international conference on Functional programming languages and computer architecture, FPCA 1989, London, UK, September 11-13, 1989, ACM, 1989, pp. 347–359. `doi:10.1145/99370.99404`.

[20] J. C. Reynolds, Types, abstraction and parametric polymorphism, in: R. E. A. Mason (Ed.), Information Processing 83, Proceedings of the IFIP 9th World Computer Congress, Paris, France, September 19-23, 1983, North-Holland/IFIP, 1983, pp. 513–523.

[21] A. Shaikhha, M. Huot, J. Smith, D. Olteanu, Functional collection programming with semi-ring dictionaries, Proc. ACM Program. Lang. 6 (OOPSLA1) (apr 2022). `doi:10.1145/3527333`.

[22] M. F. Atiyah, I. G.MacDonald, Introduction to commutative algebra., Addison-Wesley-Longman, 1969.

[23] R. Hinze, Kan extensions for program optimisation or: Art and Dan explain an old trick, in: J. Gibbons, P. Nogueira (Eds.), Mathematics of Program Construction - 11th International Conference, MPC 2012, Madrid, Spain, June 25-27, 2012. Proceedings, Vol. 7342 of Lecture Notes in Computer Science, Springer, 2012, pp. 324–362. `doi:10.1007/978-3-642-31113-0\_16`.

[24] G. Boisseau, J. Gibbons, What you needa know about Yoneda: profunctor optics and the Yoneda lemma (Functional Pearl), Proc. ACM Program. Lang. 2 (ICFP) (2018) 84:1–84:27. `doi:10.1145/3236779`.

[25] K. L. Clark, S.-Å. Tärnlund, A first order theory of data and programs, Information Processing (77) (1977).

[26] J. Hughes, A novel representation of lists and its application to the function "reverse", Inf. Process. Lett. 22 (3) (1986) 141–144. `doi:10.1016/0020-0190(86)90059-1`.

[27] A. J. Gill, J. Launchbury, S. Peyton Jones, A short cut to deforestation, in: J. Williams (Ed.), Proceedings of the conference on Functional programming languages and computer architecture, FPCA 1993, Copenhagen, Denmark, June 9-11, 1993, ACM, 1993, pp. 223–232. `doi:10.1145/165180.165214`.

[28] J. Gibbons, N. Wu, Folding domain-specific languages: deep and shallow embeddings (Functional Pearl), in: Proceedings of the 19th ACM SIGPLAN International Conference on Functional Programming, 2014, pp. 339–347.

[29] J. Karczmarczuk, Functional differentiation of computer programs, in: M. Felleisen, P. Hudak, C. Queinnec (Eds.), Proceedings of the third ACM SIGPLAN International Conference on Functional Programming (ICFP '98), Baltimore, Maryland, USA, September 27-29, 1998, ACM, 1998, pp. 195–203. `doi:10.1145/289423.289442`.

[30] R. Hinze, Functional pearl: Streams and unique fixed points, in: Proceedings of the 13th ACM SIGPLAN International Conference on Functional Programming, ICFP '08, Association for Computing Machinery, New York, NY, USA, 2008, p. 189–200. `doi:10.1145/1411204.1411232`.

[31] B. A. Pearlmutter, Fast exact multiplication by the Hessian, Neural Comput. 6 (1) (1994) 147–160. `doi:10.1162/neco.1994.6.1.147`.

[32] M. Huot, S. Staton, M. Vákár, Correctness of automatic differentiation via diffeologies and categorical gluing, in: J. Goubault-Larrecq, B. König (Eds.), Foundations of Software Science and Computation Structures - 23rd International Conference, FOSSACS 2020, Held as Part of the European Joint Conferences on Theory and Practice of Software, ETAPS 2020, Dublin, Ireland, April 25-30, 2020, Proceedings, Vol. 12077 of Lecture Notes in Computer Science, Springer, 2020, pp. 319–338. `doi:10.1007/978-3-030-45231-5\_17`.

[33] M. Huot, S. Staton, M. Vákár, Higher order automatic differentiation of higher order functions, Log. Methods Comput. Sci. 18 (1) (2022). `doi:10.46298/lmcs-18(1:41)2022`.

[34] T. Smeding, M. Vákár, Efficient dual-numbers reverse AD via well-known program transformations, `https://arxiv.org/abs/2207.03418` (2022).

[35] C. M. Elliott, Beautiful differentiation, in: G. Hutton, A. P. Tolmach (Eds.), Proceeding of the 14th ACM SIGPLAN international conference on Functional programming, ICFP 2009, Edinburgh, Scotland, UK, August 31 - September 2, 2009, ACM, 2009, pp. 191–202. `doi:10.1145/1596550.1596579`.

[36] C. Elliott, The simple essence of automatic differentiation, Proc. ACM Program. Lang. 2 (ICFP) (2018) 70:1–70:29. `doi:10.1145/3236765`.

[37] C. Elliott, Symbolic and automatic differentiation of languages, Proc. ACM Program. Lang. 5 (ICFP) (2021) 1–18. `doi:10.1145/3473583`.

[38] M. Elsman, F. Henglein, R. Kaarsgaard, M. K. Mathiesen, R. Schenck, Combinatory adjoints and differentiation, in: J. Gibbons, M. S. New (Eds.), Proceedings Ninth Workshop on Mathematically Structured Functional Programming, MSFP@ETAPS 2022, Munich, Germany, 2nd April 2022, Vol. 360 of EPTCS, 2022, pp. 1–26. `doi:10.4204/EPTCS.360.1`.

[39] F. Wang, D. Zheng, J. M. Decker, X. Wu, G. M. Essertel, T. Rompf, Demystifying differentiable programming: shift/reset the penultimate backpropagator, Proc. ACM Program. Lang. 3 (ICFP) (2019) 96:1–96:31. `doi:10.1145/3341700`.

[40] E. Kmett, `ad:` automatic differentiation, `https://github.com/ekmett/ad` (2021).

[41] B. A. Pearlmutter, J. M. Siskind, Reverse-mode AD in a functional framework: Lambda the ultimate backpropagator, ACM Trans. Program. Lang. Syst. 30 (2) (2008) 7:1–7:36. `doi:10.1145/1330017.1330018`.

[42] O. MANZYUK, B. A. PEARLMUTTER, A. A. RADUL, D. R. RUSH, J. M. SISKIND, Perturbation confusion in forward automatic differentiation of higher-order functions, Journal of Functional Programming 29 (2019) e12. `doi:10.1017/S095679681900008X`.

[43] D. D. Anderson, M. Winders, Idealization of a module, Journal of Commutative Algebra 1 (1) (2009) 3–56. `doi:10.1216/JCA-2009-1-1-3`.

[44] A. Griewank, A mathematical view of automatic differentiation, Acta Numerica 12 (2003) 321–398. `doi:10.1017/S0962492902000132`.

[45] P. Wadler, S. Blott, How to make ad-hoc polymorphism less ad hoc, in: Proceedings of the 16th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL '89, Association for Computing Machinery, New York, NY, USA, 1989, p. 60–76. `doi:10.1145/75277.75283`.
URL `https://doi.org/10.1145/75277.75283`

[46] R. Bird, Thinking Functionally with Haskell, Cambridge University Press, 2014. `doi:10.1017/CBO9781316092415`.

[47] G. Hutton, Programming in Haskell, 2nd Edition, Cambridge University Press, USA, 2016.

## Appendix A. Derivation of the Second-Order Derivative

This appendix expands on Section 6.4's claim that we can derive from the specification

$$derive_{2nd} :: (Eq\ v, Semiring\ d) \Rightarrow (v \to d) \to Expr\ v \to Dense\ v\ (Dense\ v\ d)$$
$$derive_{2nd}\ var\ e\ y\ x = eval\ var\ (tan^N\ (forwardAD_{Dense}\ Var$$
$$(tan^N\ (forwardAD_{Dense}\ Var\ e)\ y))\ x)$$

that consists of three successive appeals to *eval*, an equivalent version that is an instance of *abstractD* (which consists of a single *eval*):

$$derive'_{2nd}\ var\ x\ y\ e = tan^N\ (tan^N\ (abstractD\ gen\ e)\ x)\ y$$
$$\textbf{where}\ gen\ z = N\ (var\ z)\ (delta\ z)$$

We start the derivation from a point-free formulation of $derive_{2nd}$:

$$
\begin{aligned}
&\lambda y\ x \to\quad eval\ var \circ (\$x) \circ tan^N \circ forwardAD_{Dense}\ Var\\
&\qquad\qquad \circ (\$y) \circ tan^N \circ forwardAD_{Dense}\ Var\\
=&\ \{\text{- naturality of function application -}\}\\
&\lambda y\ x \to\quad (\$x) \circ fmap\ (eval\ var) \circ tan^N \circ forwardAD_{Dense}\ Var\\
&\qquad\qquad \circ (\$y) \circ tan^N \circ forwardAD_{Dense}\ Var\\
=&\ \{\text{- }\eta\text{-reduction -}\}\\
&\lambda y \to\quad fmap\ (eval\ var) \circ tan^N \circ forwardAD_{Dense}\ Var\\
&\qquad\qquad \circ (\$y) \circ tan^N \circ forwardAD_{Dense}\ Var\\
=&\ \{\text{- naturality of function application -}\}\\
&\lambda y \to\quad (\$y) \circ fmap\ (fmap\ (eval\ var) \circ tan^N \circ forwardAD_{Dense}\ Var)\\
&\qquad\qquad \circ tan^N \circ forwardAD_{Dense}\ Var\\
=&\ \{\text{- }\eta\text{-reduction -}\}\\
&fmap\ (fmap\ (eval\ var) \circ tan^N \circ forwardAD_{Dense}\ Var)\\
&\qquad\quad \circ tan^N \circ forwardAD_{Dense}\ Var\\
=&\ \{\text{- naturality of } tan^N; \text{ see equation (A.1) below -}\}\\
&fmap\ (tan^N \circ bimap\ (eval\ var)\ (fmap\ (eval\ var)) \circ forwardAD_{Dense}\ Var)\\
&\qquad\quad \circ tan^N \circ forwardAD_{Dense}\ Var\\
=&\ \{\text{- } eval \text{ fusion -}\}\\
&fmap\ (tan^N \circ eval\ gen') \circ tan^N \circ forwardAD_{Dense}\ Var\\
&\qquad\quad \textbf{where}\ gen'\ z = N\ (var\ z)\ (delta\ z)\\
=&\ \{\text{- } fmap \text{ fission + naturality of } tan^N \text{ -}\}\\
&fmap\ tan^N \circ tan^N \circ bimap\ (eval\ gen')\ (fmap\ (eval\ gen'))
\end{aligned}
$$

$$\circ \; forwardAD_{Dense} \; Var$$
$$\textbf{where} \; gen' \; z = N \; (var \; z) \; (delta \; z)$$
$$= \;\; \{\text{-} \; eval \; \text{fusion -}\}$$
$$fmap \; tan^N \circ tan^N \circ eval \; gen''$$
$$\textbf{where} \; gen'' \; z = N \; (N \; (var \; z) \; (delta \; z)) \; (delta \; z)$$
$$= \;\; \{\text{-} \; abstractD \; \text{definition -}\}$$
$$fmap \; tan^N \circ tan^N \circ abstractD$$
$$\textbf{where} \; gen \; z = N \; (var \; z) \; (delta \; z)$$

The steps involving the naturality of $tan^N$ do not concern the *Functor* instance of Nagata numbers, but their *Bifunctor* instance, which allows to vary both components.

> **instance** *Bifunctor* $(\cdot \ltimes \cdot)$ **where**
> $bimap \; f \; g \; (N \; d \; e) = N \; (f \; d) \; (g \; d)$

The general naturality property of $tan^N :: \forall d \; e.d \ltimes e \to e$ for this *Bifunctor* instance is:
$$g \circ tan^N = tan^N \circ bimap \; f \; g$$

However, for the specialized type $tan^N :: \forall d.d \ltimes (Dense \; v \; d) \to Dense \; v \; d$ where the two components are not independent, we also have this specialized naturality property:

$$fmap \; f \circ tan^N = tan^N \circ bimap \; f \; (fmap \; f) \tag{A.1}$$

## Appendix B. Proof of Lemma 1

*Proof:.* We prove this rule by means of the chain rule for multi-variate functions:

$$\frac{\partial f(e_1,\ldots,e_n)}{\partial x} = \sum_{i=1}^{n} \left( \left. \frac{\partial f(x_1,\ldots,x_n)}{\partial x_i} \right|_{x_1=e_1,\ldots,x_n=e_n} \times \frac{\partial e_i}{\partial x} \right)$$

Define $f$ as $f \; (y, x_1, ..., x_n) = e_2$ where $x_1, \ldots, x_n$ comprise all the free variables of $e_1$ and $e_2$. Then we have that $f(e_1, x_1, \ldots, x_n) = \textbf{let} \; y = e_1 \; \textbf{in} \; e_2$.

Using this in combination with the chain rule yields:

$$\frac{\partial(\textbf{let} \; y = e_1 \; \textbf{in} \; e_2)}{\partial x_i} \;\; = \;\; \frac{\partial f \; (e_1, x_1, ..., x_n)}{\partial x_i}$$

$$= \;\; \left. \frac{\partial f(y, x_1, \ldots, x_n)}{\partial y} \right|_{y=e_1, x_1=x_1, \ldots, x_n=x_n} \times \frac{\partial e_1}{\partial x_i}$$

$$+ \sum_{j=1}^{n} \left( \left. \frac{\partial f(y, x_1, \ldots, x_n)}{\partial x_j} \right|_{y=e_1, x_1=x_1, \ldots, x_n=x_n} \times \frac{\partial x_j}{\partial x_i} \right)$$

$$= \;\; \left. \frac{\partial e_2}{\partial y} \right|_{y=e_1} \times \frac{\partial e_1}{\partial x_i} + \left. \frac{\partial e_2}{\partial x_i} \right|_{y=e_1}$$

$$\square$$

39

## Appendix C. Haskell Libraries and Functionality

This section briefly introduces Haskell's type class feature and summarizes the key library functions that we use.

*Appendix C.1. Type Classes*

We use Haskell's type classes [45] to model the interface of algebraic structures, such as the semirings:

**class** *Semiring d* **where**
  *zero* :: *d*
  *one* :: *d*
  $(\oplus)$ :: *d → d → d*
  $(\otimes)$ :: *d → d → d*

Under the hood, the Haskell compiler creates a datatype, called a *dictionary* type, that contains all the methods of the type class.

**data** *SemiringDict d = SD { zero :: d*
                  *, one :: d*
                  *, ($\oplus$) :: d → d → d*
                  *, ($\otimes$) :: d → d → d }*

A type class instance provides the implementation of the interface for a specific type. For example,

**instance** *Semiring (Expr v)* **where**
  *zero = Zero*
  *one  = One*
  $(\oplus)$ *= Plus*
  $(\otimes)$ *= Times*

Under the hood, an instance gives rise to a dictionary value that contains these implementations.

*semiringExprDict :: SemiringDict (Expr v)*
*semiringExprDict = SD { zero = Zero*
                 *, one  = One*
                 *, ($\oplus$)  = Plus*
                 *, ($\otimes$)  = Times }*

We can use the type class methods at a specific type for which there is an instance. E.g.,

*expr :: Expr v*
*expr = one $\oplus$ one*

The compiler turns this into code that explicitly extracts the implementations from the appropriate dictionary:

$$expr :: Expr\ v$$
$$expr = d.(\oplus)\ d.one\ d.one\ \textbf{where}\ d = semiringExprDict$$

The selection of the appropriate dictionary is known as *type class resolution* and directed by the type at which the methods are used. Because every type can have at most one instance of a type class, there is always at most one appropriate dictionary.

We can also abstract over the particular type used with a type variable $d$. Then the type variable is subject to the type class constraint *Semiring* $d$. This means that $d$ can be freely chosen as long as it has a *Semiring* instance. This is known as constrained polymorphism.

$$expr :: Semiring\ d \Rightarrow d$$
$$expr = one \oplus one$$

Under the hood, the type class constraint is transformed into an explicit dictionary parameter.

$$expr :: SemiringDict\ d \rightarrow d$$
$$expr\ d = d.(\oplus)\ d.one\ d.one$$

When using such a constraint polymorphic definition at a specific type, the Haskell compiler implicitly provides the appropriate dictionary parameter. For example, it turns *expr* :: *Expr* $v$ into *expr semiringExprDict* :: *Expr* $v$.

In multi-parameter type classes, i.e., type classes with multiple parameters in their signature (e.g., *Module* $d$ $e$), we can declare *functional dependencies* between parameters to aid type class resolution. A functional dependency indicates that one type parameter is uniquely determined by the other(s). For example, our *Module* class uses such a functional dependency:

$$\textbf{class}\ (Semiring\ d, Monoid\ e) \Rightarrow Module\ d\ e\ |\ e \rightarrow d\ \textbf{where}\ ...$$

Here the semiring type $d$ can be determined from type $e$.

For a more thorough introduction to type classes and Haskell in general, we refer to Bird [46] and Hutton [47].

*Appendix C.2. Maps*

The functions below are from Haskell's *Data.Map* library[14], copied with their type signature and explanation.

- (!) :: *Ord* $k \Rightarrow Map\ k\ a \rightarrow k \rightarrow a$
  $\mathcal{O}(\log n)$. Find the value at a key. Calls *error* when the element can not be found.

---

[14]https://hackage.haskell.org/package/containers-0.4.0.0/docs/Data-Map.html

- *findWithDefault :: Ord k $\Rightarrow$ a $\rightarrow$ k $\rightarrow$ Map k a $\rightarrow$ a*
  $\mathcal{O}(\log n)$. The expression *findWithDefault def k map* returns the value at key $k$ or returns default value *def* when the key is not in the *map*.

- *fromList :: Ord k $\Rightarrow$ [(k, a)] $\rightarrow$ Map k a*
  $\mathcal{O}(n \times \log n)$. Build a map from a list of key/value pairs. If the list contains more than one value for the same key, the last value for the key is retained.

- *unionWith :: Ord k $\Rightarrow$ (a $\rightarrow$ a $\rightarrow$ a) $\rightarrow$ Map k a $\rightarrow$ Map k a $\rightarrow$ Map k a*
  $\mathcal{O}(n + m)$. Union with a combining function. The implementation uses the efficient hedge-union algorithm.

- *insertWith :: Ord k $\Rightarrow$ (a $\rightarrow$ a $\rightarrow$ a) $\rightarrow$ k $\rightarrow$ a $\rightarrow$ Map k a $\rightarrow$ Map k a*
  $\mathcal{O}(\log n)$. Insert with a function, combining new value and old value. *insertWith f key value mp* will insert the pair $(key, value)$ into *mp* if *key* does not exist in the map. If the key does exist, the function will insert the pair $(key, f \ new\_value \ old\_value)$.

*Appendix C.3. Arrays*

The functions below are from Haskell's *Data.Array.MArray* library[15], copied with their type signature and explanation.

- *newArray :: Ix i $\Rightarrow$ (i, i) $\rightarrow$ e $\rightarrow$ m (a i e)*
  Builds a new array, with every element initialised to the supplied value.

- *readArray :: (MArray a e m, Ix i) $\Rightarrow$ a i e $\rightarrow$ i $\rightarrow$ m e*
  Read an element from a mutable array

- *writeArray :: (MArray a e m, Ix i) $\Rightarrow$ a i e $\rightarrow$ i $\rightarrow$ e $\rightarrow$ m ()*
  Write an element in a mutable array

- *getAssocs :: (MArray a e m, Ix i) $\Rightarrow$ a i e $\rightarrow$ m [(i, e)]*
  Return a list of all the associations of a mutable array, in index order.

Below, we document the *ST* monad from Haskell's *Control.Monad.ST* library[16].

- **data** *ST s a*
  The ST monad allows for destructive updates, but is escapable. A computation of type *ST s a* returns a value of type *a*, and execute in "thread" *s*. It serves to keep the internal states of different invocations of runST separate from each other.

---

[15]https://hackage.haskell.org/package/array-0.5.1.1/candidate/docs/Data-Array-MArray.html
[16]https://hackage.haskell.org/package/base-4.18.0.0/docs/Control-Monad-ST.html

- $runST :: (\forall s.ST\ s\ a) \rightarrow a$
  Return the value computed by a state thread. The $\forall \cdot$ ensures that the internal state used by the $ST$ computation is inaccessible to the rest of the program.

Below, we document the $STArray$'s from Haskell's $Data.Array.ST$ library[17].

- **data** $STArray\ s\ i\ e$
  Mutable, boxed, non-strict arrays in the $ST$ monad. $s$ is the state variable argument for the $ST$ type, $i$ is the index type of the array (should be an instance of $Ix$), and $e$ is the element type of the array.

*Appendix C.4. Other*

Below, we document the $foldMap$ function from Haskell's $Data.Foldable$ library[18].

- $foldMap :: Monoid\ m \Rightarrow (a \rightarrow m) \rightarrow t\ a \rightarrow m$
  Map each element of the structure into a monoid, and combine the results with $(\oplus)$. This fold is right-associative and lazy in the accumulator.

Below, we document the $ReaderT$ monad transformer taken from Haskell's $Control.Monad.Trans.Reader$ library[19].

- **newtype** $ReaderT\ r\ m\ a = ReaderT\ \{\ runReaderT :: r \rightarrow m\ a\ \}$
  The reader monad transformer, which adds a read-only environment to the given monad $m$. The $return$ function ignores the environment, while $(\ggg\!=)$ passes the inherited environment to both subcomputations.

---

[17]https://hackage.haskell.org/package/array-0.5.5.0/docs/Data-Array-ST.html
[18]https://hackage.haskell.org/package/base-4.18.0.0/docs/Data-Foldable.html
[19]https://hackage.haskell.org/package/transformers-0.6.1.0/docs/Control-Monad-Trans-Reader.html