# Higher-Order, Higher-Order Automatic Differentiation

## *(early, incomplete draft—comments invited)*[*]

Conal Elliott

conal@conal.net

Draft of December 28, 2019 20:07

# 1   Introduction

This note picks up where Elliott (2018) left off, and I assume the reader to be familiar with that paper and have it close at hand. I am circulating this follow-on in fairly rough form for early feedback, to then evolve in to a full research paper. The main new contributions are two senses of "*higher-order* automatic differentiation":

- derivatives of higher-order functions, and

- higher-order derivatives of functions, and

The former has been addressed in a recent paper (Vytiniotis et al., 2019), but in a way I find dissatisfying for a variety of reasons described in Section 9 and discussed at length with the authors.

Begin with the category of computably differentiable functions from Elliott (2018, Section 4.1):

**newtype** $D\ a\ b = D\ (a \to b \times (a \multimap b))$

where $a \multimap b$ is the type of linear maps from $a$ to $b$. The function around which the automatic differentiation (AD) algorithm is organized simply "zips" together a function $f : a \to b$ and its derivative $\mathcal{D}\ f : a \to a \multimap b$:[1,2,3]

$$
\begin{aligned}
&\hat{\mathcal{D}} : (a \to b) \to D\ a\ b \\
&\hat{\mathcal{D}}\ f = D\ (\lambda\ a \to (f\ a, \mathcal{D}\ f\ a)) \\
&\qquad = D\ (f \mathbin{\triangle} \mathcal{D}\ f)
\end{aligned}
$$

Note that this definition is not computable, since $\mathcal{D}$ is not (Pour-El and Richards, 1978, 1983). The whole specification of AD is then simply that $\hat{\mathcal{D}}$ is a homomorphism with respect to a standard compositional vocabulary of functions, namely that of cartesian categories, plus a collection of numeric primitives like (uncurried) addition and multiplication, *sin* and *cos*, etc. An example of such a homomorphism equation is $\hat{\mathcal{D}}\ g \circ \hat{\mathcal{D}}\ f = \hat{\mathcal{D}}\ (g \circ f)$, in which the only unknown is the meaning of the LHS ($\circ$), i.e., sequential composition in the category $D$. Solving the collection of such homomorphism equations yields correct-by-construction AD.

The function $\hat{\mathcal{D}}$ is invertible, i.e., $\hat{\mathcal{D}}^{-1} \circ \hat{\mathcal{D}} = id$, where $\hat{\mathcal{D}}^{-1}$ simply drops the derivative:[4]

---

[*]For the latest version of this document and the repository location for questions, suggestions, bugs, etc, see `http://conal.net/papers/higher-order-ad`.

[1][Red bracketed remarks are notes to myself.]

[2]This paper generally uses Haskell notation but deviates slightly by using a single colon rather than double colon for type signatures. [Experimental.]

[3]The infix operators for function types ("$\to$") and linear maps ("$\multimap$") both associate to the right and have equal, very low precedence. For instance, "$a \to a \multimap b$" means $a \to (a \multimap b)$.

[4]This paper uses "*exl*" and "*exr*" to name left and right product projections (defined on cartesian categories), rather than Haskell's (function-only) "*fst*" and "*snd*".

$$\hat{\mathcal{D}}^{-1} : D \ a \ b \to (a \to b)$$
$$\hat{\mathcal{D}}^{-1} \ (D \ h) = exl \circ h$$

Indeed, $\hat{\mathcal{D}}^{-1}$ is a left inverse of $\hat{\mathcal{D}}$:

$$
\begin{aligned}
&\hat{\mathcal{D}}^{-1} \ (\hat{\mathcal{D}} \ f) \\
&= \hat{\mathcal{D}}^{-1} \ (D \ (f \vartriangle \mathcal{D} \ f)) \quad \text{-- } \hat{\mathcal{D}} \text{ definition} \\
&= exl \circ (f \vartriangle \mathcal{D} \ f) \quad\quad \text{-- } \hat{\mathcal{D}}^{-1} \text{ definition} \\
&= f \quad\quad\quad\quad\quad\quad\quad \text{-- cartesian law: } exl \circ (g \vartriangle h) = g
\end{aligned}
$$

As defined so far, $\hat{\mathcal{D}}^{-1}$ is *not* a right inverse to $\hat{\mathcal{D}}$, since the linear map portion might not be the true derivative. We will thus *restrict* the category $D$ to be the image of $\hat{\mathcal{D}}$, which is to say that $\hat{\mathcal{D}}$ is surjective, i.e., the derivative is correct.[5] This restriction guarantees that $\hat{\mathcal{D}}^{-1}$ is indeed a right inverse of $\hat{\mathcal{D}}$. Given $\hat{f} : D \ a \ b$ (with the mentioned restriction), there is an $f : a \to b$ such that $\hat{f} = \hat{\mathcal{D}} \ f$, so[6]

$$
\begin{aligned}
&\hat{\mathcal{D}} \ (\hat{\mathcal{D}}^{-1} \ \hat{f}) \\
&= \hat{\mathcal{D}} \ (\hat{\mathcal{D}}^{-1} \ (\hat{\mathcal{D}} \ f)) \quad \text{-- } \hat{f} = \hat{\mathcal{D}} \ f \\
&= \hat{\mathcal{D}} \ f \quad\quad\quad\quad\quad \text{-- } \hat{\mathcal{D}}^{-1} \circ \hat{\mathcal{D}} = id \\
&= \hat{f} \quad\quad\quad\quad\quad\quad \text{-- } \hat{\mathcal{D}} \ f = \hat{f}
\end{aligned}
$$

Thus, $\hat{\mathcal{D}} \circ \hat{\mathcal{D}}^{-1} = id$ as well.

AD is often described as coming in forward and reverse "modes". For many practical applications (including deep learning and other high-dimensional optimization problems), reverse mode is much more efficient than forward mode. As typically presented, reverse mode is also much more complicated, but this difference appears to be due only to unfortunate choices in how to understand and implement AD. Instead, a single, simple algorithm works for forward, reverse, and other modes. Reverse mode is distinguished only by using a different linear map representation resulting from a simple classic trick (Elliott, 2018).

This general AD algorithm is justified by three main theorems about differentiation:

**Theorem 1** (compose/"chain" rule).

$$\mathcal{D} \ (g \circ f) \ a = \mathcal{D} \ g \ (f \ a) \circ \mathcal{D} \ f \ a$$

**Theorem 2** (cross rule).

$$\mathcal{D} \ (f \times g) \ (a, b) = \mathcal{D} \ f \ a \times \mathcal{D} \ g \ b$$

**Theorem 3** (linear rule). *For all linear functions $f$,*

$$\mathcal{D} \ f \ a = f$$

In addition to these three theorems, we need a collection of facts about the derivatives of various mathematical operations, e.g., $\hat{\mathcal{D}} \ sin \ x = scale \ (cos \ x)$, where $scale : a \to a \multimap a$ is uncurried scalar multiplication (so $scale \ s$ is linear for all $s$).

# 2 Some Additional Properties of Differentiation

A few additional properties of differentiation will prove useful in extending Elliott (2018) to higher-order functions and higher-order derivatives.

---

[5] Haskell's type system is not expressive enough to capture this restriction by itself, so the restriction will be only implied in this draft. For more rigor, one could use a language (extension) with refinement types such as Liquid Haskell [ref] or a dependently-typed language such as Agda [ref] or Idris [ref].

[6] This reasoning hold for *any* surjective function with a left inverse.

## 2.1 Pair-Valued Domains

One half of the *curry*/*uncurry* isomorphism involves functions of pair-valued domains. The notion of partial derivatives is helpful for differentiating such functions.[7]

**Lemma 4** (proved in Section A.1). *Given a function* $f : a \times b \to c$,

$$\mathcal{D} f \, (a, b) = \mathcal{D}_l \, f \, (a, b) \triangledown \mathcal{D}_r \, f \, (a, b)$$

*where* $\mathcal{D}_l$ *and* $\mathcal{D}_r$ *construct the ("first" and "second", or "left" and "right") "partial derivatives" of* $f$ *at* $(a, b)$, *defined as follows:*

$$\mathcal{D}_l : (a \times b \to c) \to a \times b \to (a \multimap c)$$
$$\mathcal{D}_l \, f \, (a, b) = \mathcal{D} \, (f \circ (, b)) \, a$$

$$\mathcal{D}_r : (a \times b \to c) \to a \times b \to (b \multimap c)$$
$$\mathcal{D}_r \, f \, (a, b) = \mathcal{D} \, (f \circ (a, )) \, b$$

*The notation "*$(a, )$*" and "*$(b, )$*" refers to right and left "sections" of pairing:* $(, b) \, a = (a, ) \, b = (a, b)$. *Equivalently,*

$$\mathcal{D}_l \, f \, (a, b) = \mathcal{D} \, f \, (a, b) \circ inl$$
$$\mathcal{D}_r \, f \, (a, b) = \mathcal{D} \, f \, (a, b) \circ inr$$

Note also that $f \circ (a, ) = curry \, f \, a$ and $f \circ (, b) = curry' \, f \, b$, where

$$curry \, f \, a \, b = f \, (a, b)$$
$$curry' \, f \, b \, a = f \, (a, b)$$

As an example of how this decomposition of $\mathcal{D} \, f$ helps construct derivatives, suppose that $f$ is *bilinear*, which is to say that $f$ is linear in each argument, while holding the other constant. More formally *bilinearity* of $f$ means that $f \circ (a, )$ and $f \circ (b, )$ (equivalently, *curry* $f \, a$ and *curry'* $f \, b$) are both linear for all $a$ and $b$.

**Corollary 4.1.** *If* $f : a \times b \to c$ *is bilinear then*

$$\mathcal{D} \, f \, (a, b) = f \circ (, b) \triangledown f \circ (a, )$$

*Proof.*

$$\begin{aligned}
&\mathcal{D} \, f \, (a, b) \\
&= \mathcal{D}_l \, f \, (a, b) \triangledown \mathcal{D}_r \, f \, (a, b) && \text{-- Lemma 4} \\
&= \mathcal{D} \, (f \circ (, b)) \, a \triangledown \mathcal{D} \, (f \circ (a, )) \, b && \text{-- } \mathcal{D}_l \text{ and } \mathcal{D}_r \text{ definitions} \\
&= f \circ (, b) \triangledown f \circ (a, ) && \text{-- linearity}
\end{aligned}$$

$\square$

For instance, the derivative of uncurried multiplication is given by the Leibniz product rule:

$$\begin{aligned}
&\mathcal{D} \, (uncurry \, (*)) \, (a, b) \\
&= uncurry \, (*) \circ (, b) \triangledown uncurry \, (*) \circ (a, ) \\
&= (* \, b) \triangledown (a \, *) \\
&= \lambda \, (da, db) \to da * b + a * db
\end{aligned}$$

More generally, consider differentiating interacts with uncurrying:

**Corollary 4.2** (proved in Section A.2).

$$\mathcal{D} \, (uncurry \, g) \, (a, b) = at \, b \circ \mathcal{D} \, g \, a \triangledown \mathcal{D} \, (g \, a) \, b$$

---

[7]Recall that, on linear maps, $(f \triangledown g) \, (a, b) = f \, a + g \, b$, $inl \, a = (a, 0)$, and $inr \, b = (0, b)$

As a special case, let $g$ be curried multiplication:

$$\begin{aligned}
&\mathcal{D}\ (uncurry\ (*)) \\
&= at\ b \circ \mathcal{D}\ (*)\ a \triangledown \mathcal{D}\ (a*)\ b \\
&= at\ b \circ (*) \triangledown (a*) \\
&= (* b) \triangledown (a *)
\end{aligned}$$

which agrees with the calculation above.

For cartesian closure, we'll need the derivative of another function with a pair-valued domain:

$$\begin{aligned}
&eval : (a \to b) \times a \to b \\
&eval\ (f, a) = f\ a \quad \text{-- on functions}
\end{aligned}$$

(Since $eval$ is neither linear nor bilinear, Theorem 3 and Corollary 4.1 are inapplicable.) We'll need one more linear map operation, which is curried, reverse function application:[8]

$$\begin{aligned}
&at : a \to (a \to b) \multimap b \\
&at\ a\ df = df\ a
\end{aligned}$$

**Corollary 4.3** (proved in Section A.3).

$$\mathcal{D}\ eval\ (f, a) = at\ a \triangledown \mathcal{D}\ f\ a$$

## 2.2 Function-Valued Codomains

It will also be useful to calculate derivatives of functions with higher-order codomains.[9] We'll need another linear map operation, which is the indexed variant of ($\triangle$) (and a specialization of Haskell's *flip* function):

$$\begin{aligned}
&fork_F : (b \to a \multimap c) \to (a \multimap b \to c) \\
&fork_F\ h = \lambda\ da\ b \to h\ b\ da
\end{aligned}$$

**Lemma 5** (proved in Section A.4). *Given a function* $g : a \to b \to c$,

$$\mathcal{D}\ g\ a = fork_F\ (\lambda\ b \to \mathcal{D}\ (at\ b \circ g)\ a).$$

**Corollary 5.1** (proved in Section A.5).

$$\mathcal{D}\ (curry\ f)\ a = fork_F\ (\mathcal{D}_l\ f \circ (a, ))$$

# 3 Cartesian Closure, first attempt

While $D$ is a category and a *cartesian* category at that, as specified by $\hat{\mathcal{D}}$ being a cartesian functor, another question naturally arises. Can $\hat{\mathcal{D}}$ also be a cartesian *closed* functor? In other words, are there definitions of *eval*, *curry*, and *uncurry* on $D$ such that

$$\begin{aligned}
&curry\ (\hat{\mathcal{D}}\ f) = \hat{\mathcal{D}}\ (curry\ f) \\
&uncurry\ (\hat{\mathcal{D}}\ g) = \hat{\mathcal{D}}\ (uncurry\ g) \\
&eval = \hat{\mathcal{D}}\ eval
\end{aligned}$$

As usual, we'll want to solve each homomorphism equation for its single unknown, which is a categorical operation on $D$ (on the LHS).

---

[8] Linearity of $at\ a$ follows from the usual definition of addition and scaling on functions.

[9] [The previous section and this one provide "adjoint" techniques in a sense that currying is an adjunction from functions from products to functions to functions. Is there something else interesting to say here?]

## 3.1 Curry

Start with *curry*, simplifying the LHS:

$$curry\ (\hat{\mathcal{D}}\ f)$$
$$=\ curry\ (D\ (f\ \triangle\ \mathcal{D}\ f))\quad\text{-- }\hat{\mathcal{D}}\text{ definition}$$

Then the RHS:

$$\hat{\mathcal{D}}\ (curry\ f)$$

$$\begin{aligned}
&= D\ (curry\ f\ \triangle\ \mathcal{D}\ (curry\ f)) && \text{-- }\hat{\mathcal{D}}\text{ definition}\\
&= D\ (\lambda\,a \to (curry\ f\ a, \mathcal{D}\ (curry\ f)\ a)) && \text{-- }(\triangle)\text{ on functions}\\
&= D\ (\lambda\,a \to ((\lambda\,b \to f\ (a,b)), fork_F\ (\mathcal{D}_l\ f \circ (a,)))) && \text{-- Corollary 5.1}\\
&= D\ (\lambda\,a \to ((\lambda\,b \to f\ (a,b)), fork_F\ (\lambda\,b \to \mathcal{D}_l\ f\ (a,b)))) && \text{-- }(\circ)\text{ on functions}\\
&= D\ (\lambda\,a \to ((\lambda\,b \to f\ (a,b)), fork_F\ (\lambda\,b \to \mathcal{D}\ f\ (a,b) \circ inl))) && \text{-- Section A.1}
\end{aligned}$$

The last form uses $f$ and $\mathcal{D}\ f$, which can be extracted from $\hat{\mathcal{D}}\ f = D\ (f\ \triangle\ \mathcal{D}\ f)$: Thus a sufficient condition for our homomorphic specification $(curry\ (\hat{\mathcal{D}}\ f) = \hat{\mathcal{D}}\ (curry\ f))$ is

$$curry\ (D\ ff') = D\ (\lambda\,a \to ((\lambda\,b \to f\ (a,b)), fork_F\ (\lambda\,b \to f'\ (a,b) \circ inl)))$$
$$\textbf{where}\ (f,f') = fork^{-1}\ ff'$$

The $fork^{-1}$ function is half of an isomorphism that holds for all cartesian categories:

$$fork : Cartesian\ k \Rightarrow (a\ `k`\ c) \times (a\ `k`\ d) \to (a\ `k`\ (c \times d))$$
$$fork = uncurry\ (\triangle)$$

$$fork^{-1} : Cartesian\ k \Rightarrow (a\ `k`\ (c \times d)) \to (a\ `k`\ c) \times (a\ `k`\ d)$$
$$fork^{-1}\ h = (exl \circ h, exr \circ h)$$

**Lemma 6.** *The pair of functions fork and $fork^{-1}$ form an isomorphism in all cartesian categories and a linear isomorphism in the category of vector spaces and linear maps.*[10] (Proof: Exercise.)

Another such linear isomorphism can be found in cocartesian categories. The following types are specialized to biproduct categories (such as linear maps):

$$join : Cocartesian\ k \Rightarrow (a\ `k`\ c) \times (a\ `k`\ d) \to (a\ `k`\ (c \times d))$$
$$join = uncurry\ (\triangledown)$$

$$join^{-1} : Cocartesian\ k \Rightarrow (a\ `k`\ (c \times d)) \to (a\ `k`\ c) \times (a\ `k`\ d)$$
$$join^{-1}\ h = (h \circ inl, h \circ inr)$$

**Lemma 7.** *The pair of functions join and $join^{-1}$ form an isomorphism in all cocartesian categories and a linear isomorphism in the category of vector spaces and linear maps.* (Proof: Exercise.)

These two isomorphism pairs were used by Elliott (2018) to construct a correct-by-construction implementation of reverse-mode AD, by merely altering the representation of linear maps used in the simple, general AD algorithm.

Another useful operation is the *uncurried* version of the monoidal $(\times)$:

$$cross : Monoidal\ k \Rightarrow (a\ `k`\ c) \times (b\ `k`\ d) \to ((a \times b)\ `k`\ (c \times d))$$
$$cross = uncurry\ (\times)$$

**Lemma 8.** *In the category of vector spaces and linear maps, the cross function is linear.* (Proof: Exercise.)

---

[10][To do: name this category early (say "**Vec**$_s$" for a semiring $s$) and refer to by name where needed.]

Although *fork* and *fork*$^{-1}$ form an isomorphism and hence preserve information, *fork*$^{-1}$ can result in a loss of efficiency, due to computation that can be (and often is) in common to a function $f$ and its derivative $\mathcal{D} f$. Indeed, the definition of *fork*$^{-1}$ $h$ above shows that $h$ gets replicated. It's unclear how to avoid this redundancy problem in practice with currying when $D$ is used to represent computably differentiable functions. Personal experience with compiling to categories Elliott (2017) suggests that most uses of *curry* generated during translation from the $\lambda$ calculus (e.g., Haskell) are in fact transformed away at compile time using various equational CCC laws. Still, it does seem an important question to explore.

Intriguingly, curried functions can also help eliminate redundant computation suggested by uncurried counterparts functions. Given a function $g : a \to b \to c$, it is sometimes convenient to "partially apply" $g$ to an argument $u$ and then apply the resulting $g\,u : b \to c$ to many different $v : b$. In some cases, a considerable amount of work can be done based solely on $u$, saving residual work to be done for different $b$ values. In such situations, *uncurry* $g$ loses this performance advantage.

## 3.2   Uncurry

Next, let's tackle *uncurry*, whose defining homomorphism is

$$uncurry\ (\hat{\mathcal{D}}\ g) = \hat{\mathcal{D}}\ (uncurry\ g)$$

Simplify the LHS:

$$
\begin{aligned}
&uncurry\ (\hat{\mathcal{D}}\ g)\\
&= uncurry\ (D\ (g \vartriangle \mathcal{D}\ g)) \quad \text{-- } \hat{\mathcal{D} \text{ definition}
\end{aligned}
$$

Then the RHS:

$$
\begin{aligned}
&\hat{\mathcal{D}}\ (uncurry\ g)\\
&= D\ (uncurry\ g \vartriangle \mathcal{D}\ (uncurry\ g)) && \text{-- } \hat{\mathcal{D}} \text{ definition}\\
&= D\ (\lambda\,(a,b) \to (uncurry\ g\ (a,b), \mathcal{D}\ (uncurry\ g)\ (a,b))) && \text{-- } (\vartriangle) \text{ definition}\\
&= D\ (\lambda\,(a,b) \to (g\ a\ b, \mathcal{D}\ (uncurry\ g)\ (a,b))) && \text{-- } uncurry \text{ on functions}\\
&= D\ (\lambda\,(a,b) \to (g\ a\ b, at\ b \circ \mathcal{D}\ g\ a \,\triangledown\, \mathcal{D}\ (g\ a)\ b)) && \text{-- Corollary 4.2}
\end{aligned}
$$

Now we have a problem with solving the defining homomorphism above. Although we can extract $g$ and $\mathcal{D}\ g$ from $\hat{\mathcal{D}}\ g$, we cannot extract $\mathcal{D}\ (g\ a)$. Or rather we can, but not computably.

## 3.3   Eval

We don't need to work out both *uncurry* and *eval*, since each can be defined in terms of the other:

$$
\begin{aligned}
eval &= uncurry\ id\\
uncurry\ g &= eval \circ (g \times id)\\
&= eval \circ first\ g
\end{aligned}
$$

Since we got stuck on *uncurry*, let's try *eval* as well to see if we learn anything new.

The corresponding homomorphism equation has a particularly simple form:

$$eval = \hat{\mathcal{D}}\ eval$$

It might appear that we are done at the start, taking the equation to be a definition for *eval*. Recall, however, that $\hat{\mathcal{D}}$ is noncomputable, being defined via $\mathcal{D}$ (differentiation itself). Simplifying the RHS,

$$
\begin{aligned}
&\hat{\mathcal{D}}\ eval\\
&= D\ (eval \vartriangle \mathcal{D}\ eval) && \text{-- } \hat{\mathcal{D}} \text{ definition}\\
&= D\ (\lambda\,(f,a) \to (eval\ (f,a), \mathcal{D}\ eval\ (f,a))) && \text{-- } (\vartriangle) \text{ on functions}\\
&= D\ (\lambda\,(f,a) \to (f\ a, \mathcal{D}\ eval\ (f,a))) && \text{-- } eval \text{ on functions}\\
&= D\ (\lambda\,(f,a) \to (f\ a, at\ a \,\triangledown\, \mathcal{D}\ f\ a)) && \text{-- Corollary 4.3}
\end{aligned}
$$

As with uncurrying (Section 3.2), the final form is well-defined but is not a computable recipe, leaving us in a pickle. Next, let's look for some wiggle room.

# 4 Object Mapping

The choice of category-associated products and exponentials is a degree of freedom not exercised in the development of AD in Elliott (2018) (or above) and is tied closely to another such choice available in the general notion of *cartesian closed functor* in category theory. In general, a functor has two aspects:

- a mapping from arrows to arrows, and

- a mapping from objects to objects.

The functor $\hat{\mathcal{D}}$ defined (noncomputably) above implicitly chooses an *identity object mapping*, as evident in its type signature $\hat{\mathcal{D}} : (a \to b) \to D\ a\ b$. The type of $\hat{\mathcal{D}}$ plus the requirement that it be a cartesian *closed* functor implies that the object mapping aspect of $\hat{\mathcal{D}}$ is the identity. More generally, we can define an object mapping $O : Type \to Type$ for a new functor $\mathcal{D}_o$:[11]

$$\mathcal{D}_o : (a \to b) \to D\ (O\ a)\ (O\ b)$$

Each cartesian category $k$ has its own notion of categorical product $a \times_k b$ (satisfying a universality property), and similarly for cocartesian categories (with categorical products and coproducts coinciding for biproduct categories). Likewise, each cartesian *closed* category $k$ has its own notion of *exponential* objects $a \Rightarrow_k b$.

The generalized interface for cartesian closed categories with per-category exponentials is as follows:[12]

```
class Cartesian k ⇒ CartesianClosed k where
    type (⇒ₖ) : Type → Type → Type
    curry   : ((a ×ₖ b) ‘k‘ c) → (a ‘k‘ (b ⇒ₖ c))
    uncurry : (a ‘k‘ (b ⇒ₖ c)) → ((a ×ₖ b) ‘k‘ c)
    eval    : ((a ⇒ₖ b) ×ₖ a) ‘k‘ b
```

where $a \Rightarrow_k b$ is a type of "exponential objects" (first class functions/arrows) from $a$ to $b$ for the category $k$.

The property of being a closed cartesian functor requires $O$ to preserve categorical products and exponentials, i.e.,

$$O\ (a \times b) = O\ a \times_D O\ b$$
$$O\ (a \to b) = O\ a \Rightarrow_D O\ b$$

The usual notion of cartesian products are working fine, so we'll continue to choose $a \times_D b = a \times b$. While $\hat{\mathcal{D}}$ being a closed cartesian functor (CCF) from $(\to)$ to $D$ implies an noncomputable *eval* and *uncurry* (Section 3.2 and Section 3.3), our goal is to define $\Rightarrow_D$ and $\mathcal{D}_o$ such that $\mathcal{D}_o$ is a CCF with computable operations.

Consider again the homomorphic specification for *curry* (part of the CCF definition): $eval = \mathcal{D}_o\ eval$. The RHS *eval* (on functions) has type $(a \to b) \times a \to b$, while the RHS *eval* (on $D$) has type

$$
\begin{aligned}
&D\ (O\ ((a \to b) \times a))\ (O\ b) \\
=\ &D\ (O\ (a \to b) \times O\ a)\ (O\ b) \\
=\ &D\ ((O\ a \Rightarrow_D O\ b) \times O\ a)\ (O\ b)
\end{aligned}
$$

The difficulty with our attempt at *eval* in Section 3.3 was that we were given a function $f$, but we also needed its derivative $\mathcal{D}\ f$. Similarly, with *uncurry* in Section 3.2, we were given $g : a \to b \to c$, and we needed not only $g\ a$ but also $\mathcal{D}\ (g\ a)$. In both cases the exponential object was a function, but we also needed its (computable) derivative.

This analysis suggests that we include a derivative in the exponential object, simply by choosing $\Rightarrow_D$ to be $D$ itself. Additionally, map scalars to themselves and cartesian products to cartesian products:

$$
\begin{aligned}
O\ \mathbb{R} &= \mathbb{R} \\
O\ (a \times b) &= O\ a \times_D O\ b = O\ a \times O\ b \\
O\ (a \to b) &= O\ a \Rightarrow_D O\ b = D\ (O\ a)\ (O\ b)
\end{aligned}
$$

---

[11][Experiment with different notation for $O\ a$, e.g., "$\bar{a}$".]

[12]These operations support higher-order programming and arise during translation from a typed lambda calculus (e.g., Haskell) to categorical vocabulary (Elliott, 2017).

We will need to convert between $a$ and $O\ a$, which we can do with a family of *linear isomorphisms*[13] indexed by $a$:[14]

> **class** $HasO\ t$ **where**
>> **type** $O\ t$
>> $o\quad: t \to O\ t$
>> $o^{-1}: O\ t \to t$

For scalar types $a$ and the unit type, $O\ a = a$, the isomorphism is trivial:

> **instance** $HasO\ \mathbb{R}$ **where**
>> **type** $O\ \mathbb{R} = \mathbb{R}$
>> $o\quad = id$
>> $o^{-1} = id$

> **instance** $HasO\ ()$ **where**
>> **type** $O\ () = ()$
>> $o\quad = id$
>> $o^{-1} = id$

For products, convert components independently:[15]

> **instance** $(HasO\ a, HasO\ b) \Rightarrow HasO\ (a \times b)$ **where**
>> **type** $O\ (a \times b) = O\ a \times O\ b$
>> $o\quad = o\quad \times o$
>> $o^{-1} = o^{-1} \times o^{-1}$

The new functor $\mathcal{D}_o$ converts its given $a \to b$ to $O\ a \to O\ b$ and then applies the $\hat{\mathcal{D}}$ functor:[16]

$$(\Rightarrow) : (p' \to p) \to (q \to q') \to ((p \to q) \to (p' \to q'))$$
$$f \Rightarrow h = \lambda\, g \to h \circ g \circ f$$

$$wrap_o : (a \to b) \to (O\ a \to O\ b)$$
$$wrap_o = o^{-1} \Rightarrow o$$

$$wrap_o^{-1} : (O\ a \to O\ b) \to (a \to b)$$
$$wrap_o^{-1} = o \Rightarrow o^{-1}$$

$$\mathcal{D}_o : (a \to b) \to D\ (O\ a)\ (O\ b)$$
$$\mathcal{D}_o = \hat{\mathcal{D}} \circ wrap_o$$

$$\mathcal{D}_o^{-1} : D\ (O\ a)\ (O\ b) \to (a \to b)$$
$$\mathcal{D}_o^{-1} = wrap_o^{-1} \circ \hat{\mathcal{D}}^{-1}$$

**Lemma 9** (proved in Section A.6). $wrap_o$ and $wrap_o^{-1}$ *form a linear isomorphism.*

**Lemma 10** (proved in Section A.7). $\mathcal{D}_o$ and $\mathcal{D}_o^{-1}$ *form a linear isomorphism.*

**Lemma 11** (proved in Section A.8). $wrap_o$ *is a cartesian functor.*

The cartesian category operations already defined on $D$ (Elliott, 2018) are solutions to homomorphism equations saying that $\hat{\mathcal{D}}$ is a cartesian functor. Thanks to the simple, regular structure of $o$ and $o^{-1}$,

**Theorem 12.** $\mathcal{D}_o$ *is a cartesian functor.*

---

[13]The implicit requirements for all $HasO$ instances are thus that $o \circ o^{-1} = id$, $o^{-1} \circ o = id$, and *to* and $o^{-1}$ are linear.
[14][It may be more elegant to combine the functions $o$ and $o^{-1}$ into a single *isomorphism*.]
[15]Recall that $(f \times g)\ (a, b) = (f\ a, g\ b)$, so $o\ (a, b) = (o \times o)\ (a, b) = (o\ a, o\ b)$, and similarly for $o^{-1}$.
[16][Consider dropping the $(\Rightarrow)$ definition and uses here.]

Proof: $\hat{\mathcal{D}}$ is a cartesian functor (Elliott, 2018), as is $wrap_o$ (Lemma 11), so $\mathcal{D}_o = \hat{\mathcal{D}} \circ wrap_o$ is also.

What about exponentials and cartesian *closure*? As mentioned above, $O\ (a \to b) = O\ a \Rightarrow_D O\ b = D\ (O\ a)\ (O\ b)$, which suggests using $\mathcal{D}_o$ and $\mathcal{D}_o^{-1}$ for $o$ and $o^{-1}$:

> **instance** $(HasO\ a, HasO\ b) \Rightarrow HasO\ (a \to b)$ **where**
> **type** $O\ (a \to b) = D\ (O\ a)\ (O\ b)$
> $o\ \ = \mathcal{D}_o$
> $o^{-1} = \mathcal{D}_o^{-1}$

A useful consequence:

**Lemma 13** (proved in Section A.9)**.**

$$wrap_o\ (curry\ f) = \hat{\mathcal{D}} \circ curry\ (wrap_o\ f)$$

**Corollary 13.1.**

$$curry\ (wrap_o\ f) = \hat{\mathcal{D}}^{-1} \circ wrap_o\ (curry\ f)$$

*Proof.* Left-compose $\hat{\mathcal{D}}^{-1}$ with both sides of Lemma 13; then simplify and reverse the resulting equation. □

Let's now try to solve the CCF equations for $\mathcal{D}_o$. This time begin with *eval*:

**Lemma 14** (proved in Section A.10)**.** *With the following (effective) definition of eval on $D$, eval $= \mathcal{D}_o$ eval:*

$$eval = D\ (\lambda\ (D\ h, a) \to \mathbf{let}\ (b, f') = h\ a\ \mathbf{in}\ (b, at\ a \circ \hat{\mathcal{D}}^{-1}\ \triangledown\ f'))$$

For *uncurry*, use the standard definition $uncurry\ g = eval \circ first\ g$.

The definition of *curry* in Section 3.1 worked fine, but we'll need to check again, as we did with the cartesian category operations (Theorem 12). The homomorphism equation is $curry\ (\mathcal{D}_o\ f) = \mathcal{D}_o\ (curry\ f)$, to be solved for the unknown LHS $curry$ (on $D$), with $f : a \times b \to c$. First let $f_o = wrap_o\ f$. Simplify the LHS:

> $curry\ (\mathcal{D}_o\ f)$
> $= curry\ (\hat{\mathcal{D}}\ (wrap_o\ f))$    -- $\mathcal{D}_o$ definition
> $= curry\ (\hat{\mathcal{D}}\ f_o)$    -- $f_o$ definition
> $= curry\ (D\ (f_o \vartriangle \mathcal{D}\ f_o))$    -- $\hat{\mathcal{D}}$ definition

Then the RHS:[17]

**Lemma 15** (proved in Section A.11)**.**

> $\mathcal{D}_o\ (curry\ f) =$
>      $D\ (\lambda\ a \to (D\ (\lambda\ b \to (f_o\ (a, b), \mathcal{D}_r\ f_o\ (a, b)))$
>             $, \lambda\ da \to D\ (\lambda\ b \to (\mathcal{D}_l\ f_o\ (a, b)\ da, at\ da \circ \mathcal{D}_r\ (\mathcal{D}_l\ f_o)\ (a, b)))))$

*where $f_o = wrap_o\ f$.*

The RHS uses $f_o\ (a, b)$ and $\mathcal{D}\ f_o\ (a, b)$ (via its components $\mathcal{D}_l\ f_o\ (a, b)$ and $\mathcal{D}_r\ f_o\ (a, b)$), but it also uses a *second* partial derivative $\mathcal{D}_r\ (\mathcal{D}_l\ f_o)\ (a, b)$, which is not available from the *curry* argument $D\ (f_o \vartriangle \mathcal{D}\ f_o)$.

---

[17][State, prove, and use a lemma about $\hat{\mathcal{D}}\ (g \circ f)\ a$ for linear $g$ and another for linear $f$. Maybe also $\mathcal{D}_o\ (g \circ f)\ a$ for linear $g$ or $f$.]

# 5 Where Are We?

Let's now reflect on what we've learned so far:

- The cartesian functor (CF) $\hat{\mathcal{D}} : (a \to b) \to D\ a\ b$ also forms a cartesian *closed* functor (CCF) with suitable definitions of *curry*, *uncurry*, and *eval*, but not computably (Section 3). More specifically, *curry* is computable, but *uncurry* and *eval* are not, since they need to synthesize derivatives of regular computable functions.

- General categorical functors can remap objects (here, types) as well as morphisms (here, functions). Exploiting this degree of freedom, define $\mathcal{D}_o : (a \to b) \to D\ (O\ a)\ (O\ b)$, where $O : Type \to Type$ replaces regular functions with computably differentiable functions, i.e., $O\ (u \to v) = D\ (O\ u)\ (O\ v)$. This new function is defined in terms of the old one, $\mathcal{D}_o = \hat{\mathcal{D}} \circ wrap_o$, and indeed $\mathcal{D}_o$ is a CF as well. In the absence of higher-order functions, $O$ is the identity mapping, and $\mathcal{D}_o$ coincides with $\hat{\mathcal{D}}$.

- Computably satisfying the required homomorphism properties of $\mathcal{D}_o$ for *uncurry* and *eval* becomes easy, since the operations are *given* the required derivatives rather than having to synthesize them. Unfortunately, now *curry* becomes noncomputable because it has to synthesize partial *second* derivatives.

# 6 Higher-Order Derivatives

Where can we go from here? An obvious next step is to add second order derivatives to the representation of computably differentiable functions. It seem likely, however, that the CCF specification would reveal that *curry* needs at least third order derivatives, and so on. In other words, differentiation of higher-order functions requires all higher-order derivatives of functions.

In order to construct higher-order derivatives, it will help to examine the linearity properties of our familiar categorical vocabulary, which turns out to be mostly linear with just a bit of bilinearity. As noted in Elliott (2018), the categorical operation *id*; the cartesian operations *exl*, *exr*, *dup*; and the cocartesian operations *inl*, *inr*, and *jam* are all linear. Lemmas 6 and 7 have already noted that the functions *fork* and *join* (uncurried versions of ($\triangle$) and ($\triangledown$) defined in Section 3.1) are linear (as well as isomorphisms). Next, let *comp* be uncurried composition:[18]

$$comp : Category\ k \Rightarrow (b\ `k`\ c) \times (a\ `k`\ b) \to (a\ `k`\ c)$$
$$comp = uncurry\ (\circ)$$

**Lemma 16** (proved in Section A.12). *On linear maps, comp is bilinear.*

**Lemma 17** (proved in Section A.13). *Given any bilinear function h:*

  a. *curry h a is linear for all a.*

  b. *curry$'$ h b is linear for all b.*

  c. *curry h and curry$'$ h are linear.*

  d. $\mathcal{D}\ h$ *is linear.*

**Corollary 17.1.** *On linear maps,*

  a. $(g \circ)$ *is linear for all g.*

  b. $(\circ f)$ *is linear for all f.*

  c. $(\circ)$ *and flip* $(\circ)$ *are linear.*

  d. $\mathcal{D}\ comp$ *is linear.*

These properties will help re-express Theorems 1 and 2 and related facts in a form more amenable to constructing higher derivatives:

---

[18][Maybe define *comp* only for linear maps.]

**Lemma 18** (proved in Section A.14)**.**

   a. $\mathcal{D}\,(g \circ f) = comp \circ (\mathcal{D}\,g \circ f \,\vartriangle\, \mathcal{D}\,f)$.

   b. $\mathcal{D}\,(f \times g) = cross \circ (\mathcal{D}\,f \times \mathcal{D}\,g)$.

   c. $\mathcal{D}\,(f \,\vartriangle\, g) = fork \circ (\mathcal{D}\,f \,\vartriangle\, \mathcal{D}\,g)$.

   d. *For a* linear *function* $f$, $\mathcal{D}\,f = const\,f$.

   e. *For any function* $f : a \times b \to c$, $\mathcal{D}\,f = join \circ (\mathcal{D}_l\,f \times \mathcal{D}_r\,f)$.

   f. *For a* bilinear *function* $f : a \times b \to c$, $\mathcal{D}\,f = join \circ (curry'\,f \times curry\,f) \circ swap$.

   g. *On linear maps,* $\mathcal{D}\,comp = join \circ (flip\,(\circ) \times (\circ)) \circ swap$.

Let us now consider the task of constructing *all* orders of derivatives. The $D$ category encapsulates a function $f$ and its first derivative, i.e., the zeroth and first derivatives of $f$, which we might write as "$\hat{\mathcal{D}}\,f = \mathcal{D}^0\,f \,\vartriangle\, \mathcal{D}^1\,f$". Our new category will encapsulate *all* derivatives of $f$, i.e.,[19]

$$\mathcal{D}^*\,f = \mathcal{D}^0\,f \,\vartriangle\, \mathcal{D}^1\,f \,\vartriangle\, \mathcal{D}^2\,f \,\vartriangle\, \cdots$$

where

$$\begin{aligned} \mathcal{D}^0 \quad f &= f \\ \mathcal{D}^{n+1}\,f &= \mathcal{D}^n\,(\mathcal{D}\,f) \end{aligned}$$

Then

$$\begin{aligned} &\mathcal{D}^*\,f \\ &= \mathcal{D}^0\,f \,\vartriangle\, \mathcal{D}^1\,f \,\vartriangle\, \mathcal{D}^2\,f \,\vartriangle\, \mathcal{D}^3\,f \,\vartriangle\, \cdots \\ &= f \,\vartriangle\, \mathcal{D}^1\,f \,\vartriangle\, \mathcal{D}^2\,f \,\vartriangle\, \mathcal{D}^3\,f \,\vartriangle\, \cdots \\ &= f \,\vartriangle\, \mathcal{D}^0\,(\mathcal{D}\,f) \,\vartriangle\, \mathcal{D}^1\,(\mathcal{D}\,f) \,\vartriangle\, \mathcal{D}^2\,(\mathcal{D}\,f) \,\vartriangle\, \cdots \\ &= f \,\vartriangle\, \mathcal{D}^*\,(\mathcal{D}\,f) \end{aligned}$$

which we can take as a recursive definition of $\mathcal{D}^*$. Define a corresponding type of infinitely differentiable functions:[20]

   **type** $D^*\ a\ b = a \to T\ a\ b$

   **type** $T\ a\ b = b \times T\ a\ (a \multimap b)$

   $\mathcal{D}^* : (a \to b) \to D^*\ a\ b$
   $\mathcal{D}^*\,f = f \,\vartriangle\, \mathcal{D}^*\,(\mathcal{D}\,f)$

We will want to find cartesian category operations for $D^*$ such that $\mathcal{D}^*$ is a cartesian functor (CF), which will be coinductively assumed at several points below.

Start with the constant-zero function[21]: $zero : a \to b$:

$$\begin{aligned} &\mathcal{D}^*\,zero \\ &= zero \,\vartriangle\, \mathcal{D}^*\,(\mathcal{D}\,zero) &&\text{-- } \mathcal{D}^*\ \text{definition} \\ &= zero \,\vartriangle\, \mathcal{D}^*\,zero &&\text{-- } \mathcal{D}\,zero = const\,zero = zero \\ &= zero \,\vartriangle\, zero &&\text{-- coinduction} \\ &= zero &&\text{-- Zero on pairs} \end{aligned}$$

Then constant functions more generally:

---

[19]Take $\vartriangle$ to be *right*-associative.

[20]For notational simplicity, we'll drop the **newtype** isomorphisms.

[21]As usual, types are restricted to vector spaces over a common field, which we can take to be $\mathbb{R}$

$$
\begin{aligned}
&\mathcal{D}^*\,(const\ b)\\
&= const\ b \vartriangle \mathcal{D}^*\,(\mathcal{D}\,(const\ b)) \quad \text{-- } \mathcal{D}^* \text{ definition}\\
&= const\ b \vartriangle \mathcal{D}^*\ zero \quad\qquad\qquad \text{-- } \mathcal{D}\,(const\ b) = zero\\
&= const\ b \vartriangle zero \qquad\qquad\qquad\ \text{-- above}
\end{aligned}
$$

Next, linear functions $f$:

$$
\begin{aligned}
&\mathcal{D}^*\,f\\
&= f \vartriangle \mathcal{D}^*\,(\mathcal{D}\,f) \qquad\quad \text{-- } \mathcal{D}^* \text{ definition}\\
&= f \vartriangle \mathcal{D}^*\,(const\ f) \quad\ \text{-- } f \text{ linearity}\\
&= f \vartriangle const\ f \vartriangle zero \quad \text{-- above}
\end{aligned}
$$

We will have several uses of this formula, so name it:

$$
\begin{aligned}
&linear : (a \multimap b) \to D^*\ a\ b\\
&linear\ f = f \vartriangle const\ f \vartriangle zero
\end{aligned}
$$

For instance, the following definitions of $id$, $exl$ and $exr$ satisfy the associated homomorphism (cartesian functor) properties:

$$
\begin{aligned}
id\ &= linear\ id\\
exl\ &= linear\ exl\\
exr\ &= linear\ exr
\end{aligned}
$$

Next, *bilinear* functions $g$:

$$
\begin{aligned}
&\mathcal{D}^*\,g\\
&= g \vartriangle \mathcal{D}^*\,(\mathcal{D}\,g) \qquad\qquad\qquad\qquad\qquad\quad \text{-- } \mathcal{D}^* \text{ definition}\\
&= g \vartriangle linear\,(\mathcal{D}\,g) \qquad\qquad\qquad\qquad\qquad\ \text{-- derivative of bilinear is linear}\\
&= g \vartriangle linear\,(join \circ (curry'\ g \times curry\ g) \circ swap) \quad \text{-- Lemma 18f}
\end{aligned}
$$

Specialize to uncurried linear map composition:

$$
\begin{aligned}
&\mathcal{D}^*\,comp\\
&= comp \vartriangle linear\,(join \circ (curry'\ comp \times curry\ comp) \circ swap) \quad \text{-- above}\\
&= comp \vartriangle linear\,(join \circ (flip\ (\circ) \times (\circ)) \circ swap) \qquad\qquad \text{-- } comp \text{ definition}
\end{aligned}
$$

Name $\mathcal{D}^*\,comp$ for future use:

$$
\begin{aligned}
&comp' : D^*\,((b \multimap c) \times (a \multimap b))\,(a \multimap c)\\
&comp' = comp \vartriangle linear\,(join \circ (flip\ (\circ) \times (\circ)) \circ swap)
\end{aligned}
$$

Then sequential compositions:

$$
\begin{aligned}
&\mathcal{D}^*\,(g \circ f)\\
&= g \circ f \vartriangle \mathcal{D}^*\,(\mathcal{D}\,(g \circ f)) \qquad\qquad\qquad\qquad \text{-- } \mathcal{D}^* \text{ definition}\\
&= g \circ f \vartriangle \mathcal{D}^*\,(comp \circ (\mathcal{D}\,g \circ f \vartriangle \mathcal{D}\,f)) \qquad\quad \text{-- Lemma 18a}\\
&= g \circ f \vartriangle \mathcal{D}^*\,comp \circ (\mathcal{D}^*\,(\mathcal{D}\,g) \circ \mathcal{D}^*\,f \vartriangle \mathcal{D}^*\,(\mathcal{D}\,f)) \quad \text{-- coinduction}\\
&= g \circ f \vartriangle comp' \circ (\mathcal{D}^*\,(\mathcal{D}\,g) \circ \mathcal{D}^*\,f \vartriangle \mathcal{D}^*\,(\mathcal{D}\,f)) \qquad \text{-- above}
\end{aligned}
$$

Note that all of the components here ($g$, $f$, $\mathcal{D}^*\,(\mathcal{D}\,g)$, $\mathcal{D}^*\,f$, and $\mathcal{D}^*\,(\mathcal{D}\,f)$) are available in $\mathcal{D}^*\,g$ and $\mathcal{D}^*\,f$, so we have a computable recipe for $(\circ)$ on $D^*$. [To do: fill in the details.]

Finally, $f \vartriangle g$:

$$
\begin{aligned}
&\mathcal{D}^*\,(f \vartriangle g)\\
&= (f \vartriangle g) \vartriangle \mathcal{D}^*\,(\mathcal{D}\,(f \vartriangle g)) \qquad\qquad\qquad \text{-- } \mathcal{D}^* \text{ definition}\\
&= (f \vartriangle g) \vartriangle \mathcal{D}^*\,(fork \circ (\mathcal{D}\,f \vartriangle \mathcal{D}\,g)) \qquad\ \text{-- Lemma 18c}\\
&= (f \vartriangle g) \vartriangle \mathcal{D}^*\,fork \circ (\mathcal{D}^*\,(\mathcal{D}\,f) \vartriangle \mathcal{D}^*\,(\mathcal{D}\,g)) \qquad \text{-- coinduction}\\
&= (f \vartriangle g) \vartriangle linear\ fork \circ (\mathcal{D}^*\,(\mathcal{D}\,f) \vartriangle \mathcal{D}^*\,(\mathcal{D}\,g)) \quad \text{-- } fork \text{ linearity (Lemma 6)}
\end{aligned}
$$

Again, the components here ($f$, $g$, $\mathcal{D}^*\,(\mathcal{D}\,f)$, and $\mathcal{D}^*\,(\mathcal{D}\,g)$) are all available from $\mathcal{D}^*\,f$ and $\mathcal{D}^*\,g$, so we have a computable recipe for $(\vartriangle)$ on $D^*$. [To do: fill in the details.]

# 7 Avoiding redundant computation

The $\hat{\mathcal{D}}$ functor was carefully chosen to enable elimination of redundant computation between a function and its derivative. The potential for redundancy is apparent in the chain rule (Theorem 1):

$$\mathcal{D} \ (g \circ f) \ a = \mathcal{D} \ g \ (f \ a) \circ \mathcal{D} \ f \ a$$

This theorem reveals that computation of $(g \circ f) \ a$ and $\mathcal{D} \ (g \circ f) \ a$ at both involve computing $f \ a$. Since sequential composition is a very commonly used building block of computations, it is thus typical for functions and their derivatives to involve common work. This fact motivates the choice $\hat{\mathcal{D}} \ f = f \vartriangle \mathcal{D} \ f$ over $\mathcal{D}_0^+ \ f = (f, \mathcal{D} \ f)$ (Elliott, 2018, Section 3.1). While both options can give rise to compositional (functorial) AD, $\mathcal{D}_0^+$ precludes sharing of work, while $\hat{\mathcal{D}}$ enables such sharing, with just a bit of care:

$$D \ g \circ D \ f = D \ (\lambda \ a \to \textbf{let} \ \{(b, f') = f \ a; (c, g') = g \ b\} \ \textbf{in} \ (c, g' \circ f'))$$

We can calculate this definition in a categorical/pointfree form using Lemma 18a:[22]

$$
\begin{aligned}
&\hat{\mathcal{D}} \ (g \circ f) \\
&= D \ (g \circ f \vartriangle \mathcal{D} \ (g \circ f)) && \text{-- } \hat{\mathcal{D}} \text{ definition} \\
&= D \ (g \circ f \vartriangle comp \circ (\mathcal{D} \ g \circ f \vartriangle \mathcal{D} \ f)) && \text{-- Lemma 18a} \\
&= D \ (second \ comp \circ (g \circ f \vartriangle (\mathcal{D} \ g \circ f \vartriangle \mathcal{D} \ f))) && \text{-- Lemma 19c below} \\
&= D \ (second \ comp \circ assoc_R \circ ((g \circ f \vartriangle \mathcal{D} \ g \circ f) \vartriangle \mathcal{D} \ f)) && \text{-- [justify this step]} \\
&= D \ (second \ comp \circ assoc_R \circ ((g \vartriangle \mathcal{D} \ g) \circ f \vartriangle \mathcal{D} \ f)) && \text{-- Gibbons (2002, Section 1.5.1).} \\
&= D \ (second \ comp \circ assoc_R \circ (\hat{\mathcal{D}} \ g \circ f \vartriangle \mathcal{D} \ f)) && \text{-- } \hat{\mathcal{D}} \text{ definition} \\
&= D \ (second \ comp \circ assoc_R \circ first \ (\hat{\mathcal{D}} \ g) \circ (f \vartriangle \mathcal{D} \ f)) && \text{-- Lemma 19b below} \\
&= D \ (second \ comp \circ assoc_R \circ first \ (\hat{\mathcal{D}} \ g) \circ \hat{\mathcal{D}} \ f) && \text{-- } \hat{\mathcal{D}} \text{ definition}
\end{aligned}
$$

We can thus define

$$D \ \hat{g} \circ D \ \hat{f} = D \ (second \ comp \circ assoc_R \circ first \ \hat{g} \circ \hat{f})$$

with the consequence that $\hat{\mathcal{D}} \ g \circ \hat{\mathcal{D}} \ f = \hat{\mathcal{D}} \ (g \circ f)$. In this form, $\hat{f}$ and $\hat{g}$ each appear once, so as long as $D \ \hat{f}$ and $D \ \hat{g}$ are nonredundant, $D \ \hat{g} \circ D \ \hat{f}$ will be nonredundant as well. Inlining the definitions of $comp$ and of $second$, $assoc_R$, and $first$ for functions and then simplifying yields the pointful definition above.

**Lemma 19.** *The following properties hold for* $(\vartriangle)$:

a. $(h \times k) \circ (f \vartriangle g) = h \circ f \vartriangle k \circ g$

b. $first \ h \circ (f \vartriangle g) = h \circ f \vartriangle g$

c. $second \ k \circ (f \vartriangle g) = f \vartriangle k \circ g$

Proof: For $a$, see Gibbons (2002, Section 1.5.1). Then $b$ and $c$ follow as corollaries from the definitions $first \ h = h \times id$ and $second \ k = id \times k$.

---

[22]The $assoc_R$ operation in monoidal categories is defined for functions as $assoc_R \ ((a, b), c) = (a, (b, c))$.

# 8 What's Next?

<span style="color:red">[Yet to come:</span>

- <span style="color:red">Avoid redundant computation in $D^*$. Doing so is fairly easy in $D$ (zeroth and first derivatives), but I don't yet see how in $D^*$ (all derivatives).</span>

- <span style="color:red">Spell out the *Category* and *Cartesian* instances that result from solving the cartesian functor equations as in Section 6.</span>

- <span style="color:red">Cartesian *closure* (*curry* and *eval*/*uncurry*) for $D^*$, exploiting higher-order derivatives.</span>

- <span style="color:red">Variation of $\mathcal{D}^*$: $\mathcal{D}_* \, f = f \mathbin{\triangle} \mathcal{D} \, (\mathcal{D}_* \, f)$.</span>

<span style="color:red">]</span>

# 9 Related Work

The most closely related work I'm aware of is by Vytiniotis et al. (2019)[23], who also define an algorithm around the language of cartesian closed categories. There appear to be some significant shortcomings, however, at least when considered as an extension to Elliott (2018):

- Although the work is referred to as "differentiable programming", it appears to lack a specification and proof that match this claim, i.e., one defined by the mathematical operation of differentiation. As such, it's unclear to me whether the algorithm is about differentiation or something else. In contrast, the specification at the center of Elliott (2018) (and the extensions described above) is just (Fréchet) differentiation itself, combined with the original function as needed by the chain rule, or rather the requirement that the function-with-derivative satisfies a standard collection of homomorphism properties. Correctness of the algorithm was defined as faithfulness to this simple specification, and the algorithm is systematically derived from this specification and hence is correct by construction.

- Functions are already well-defined as a vector space, and thus linear maps (including derivatives) are as well, but the authors chose a different notion. They write

  > [...] what should be the tangent space of a function type? Perhaps surprisingly, a function type itself is not the right answer. We provide two possible implementations for function tangents and differentiable currying, and explain the tradeoffs.

  There is no explanation, however, of what makes their answers "right" and the unsurprising answer wrong. It is unclear what it could possibly mean for their answer to be right, since the usual notion of derivative of a function $f : a \to b$ between vector spaces has type $a \to a \multimap b$ for all vector spaces $a$ and $b$, *including function types*. This observation seems to contradict the claim that the tangent space for a function types is not a function type.

- The algorithm presented is limited to reverse mode rather than a general AD algorithm as in Elliott (2018) and the work described above.

Another related paper is Brunel et al. (2019). The authors write (in Section 1)

> However, Elliot's approach is still restricted to first-order programs (i.e., computational graphs): as far as we understand, the functor D is cartesian but not cartesian closed, so the higher-order primitives ($\lambda$-abstraction and application) lack a satisfactory treatment. This is implicit in Sect. 4.4 of Elliott (2018), where the author states that he only uses biproduct categories: it is well-known that non-trivial cartesian closed biproduct categories do not exist.

The confusion here—which was mistakenly encouraged by Elliott (2018)—is the idea that the category of differentiable functions itself is (or need be) a biproduct category. Rather, all that was needed is that the various representations of *linear maps* (derivatives) are biproduct categories. This requirement is easily satisfied by construction, since these representations are all calculated from their denotation (linear functions, itself a biproduct category) via simple cocartesian functors.

---

[23]I am in the middle of an in-depth conversation with authors.

# References

Aloïs Brunel, Damiano Mazza, and Michele Pagani. Backpropagation in the simply typed lambda-calculus with linear negation. *CoRR*, abs/1909.13768, 2019.

Conal Elliott. Compiling to categories. In *Proceedings of the ACM on Programming Languages (ICFP)*, 2017.

Conal Elliott. The essence of automatic differentiation. In *Proceedings of the ACM on Programming Languages (ICFP)*, 2018.

Jeremy Gibbons. Calculating functional programs. In *Algebraic and Coalgebraic Methods in the Mathematics of Program Construction*, volume 2297 of *Lecture Notes in Computer Science*. Springer-Verlag, 2002.

Marian Boykan Pour-El and Ian Richards. Differentiability properties of computable functions—A summary. *Acta Cybernetica*, 4(1):123–125, 1978.

Marian Boykan Pour-El and Ian Richards. Computability and noncomputability in classical analysis. *Transactions of the American Mathematical Society*, 275(2):539–560, 1983.

Dimitrios Vytiniotis, Dan Belov, Richard Wei, Gordon Plotkin, and Martin Abadi. The differentiable curry. October 2019. To appear in the Program Transformations workshop at NeurIPS 2019.

# A   Proofs

## A.1   Lemma 4

Suppose we have a function $f : a \times b \to c$, and we want to compute its derivative at a point in its (pair-valued) domain. Because linear maps (derivatives) form a cocartesian category,[24]

$$\mathcal{D} f\ (a, b) = \mathcal{D} f\ (a, b) \circ inl \triangledown \mathcal{D} f\ (a, b) \circ inr$$

Noting that (for linear maps) $inl\ da = (da, 0)$ and $inr\ db = (0, db)$, we can see that the "partial derivatives" ($\mathcal{D} f\ (a, b) \circ inl$ and $\mathcal{D} f\ (a, b) \circ inr$) allow only one half of a pair to change.

Next, note that $\mathcal{D} f\ (a, b) \circ inl = \mathcal{D}\ (f \circ (, b))\ a$, by the following equational reasoning:

$$
\begin{aligned}
&\mathcal{D}\ (f \circ (, b))\ a \\
&= \mathcal{D} f\ ((, b)\ a) \circ \mathcal{D}\ (, b)\ a && \text{-- chain rule (Theorem 1)} \\
&= \mathcal{D} f\ (a, b) \circ \mathcal{D}\ (, b)\ a && \text{-- } (, b) \text{ definition} \\
&= \mathcal{D} f\ (a, b) \circ \mathcal{D}\ (inl + const\ (0, b))\ a && \text{-- } inl \text{ on functions, and meaning of } (, b) \\
&= \mathcal{D} f\ (a, b) \circ (\mathcal{D}\ inl\ a + \mathcal{D}\ (const\ (0, b))\ a) && \text{-- linearity of } (+) \\
&= \mathcal{D} f\ (a, b) \circ \mathcal{D}\ inl\ a && \text{-- } \mathcal{D}\ (const\ z)\ a = 0 \\
&= \mathcal{D} f\ (a, b) \circ inl && \text{-- linearity of } inl; \text{ Theorem 3}
\end{aligned}
$$

Likewise, $\mathcal{D} f\ (a, b) \circ inr = \mathcal{D}\ (f \circ (a, ))\ b$.

## A.2   Corollary 4.2

$$
\begin{aligned}
&\mathcal{D}\ (uncurry\ g)\ (a, b) \\
&= \mathcal{D}_l\ (uncurry\ g)\ (a, b) \triangledown \mathcal{D}_r\ (uncurry\ g)\ (a, b) && \text{-- Lemma 4} \\
&= \mathcal{D}\ (uncurry\ g \circ (, b))\ a \triangledown \mathcal{D}\ (uncurry\ g \circ (a, ))\ b && \text{-- } \mathcal{D}_l \text{ and } \mathcal{D}_r \text{ definitions} \\
&= \mathcal{D}\ (\lambda\, a' \to uncurry\ g\ (a', b))\ a \triangledown && \text{-- } \eta \text{ expansion and simplification} \\
&\quad \mathcal{D}\ (\lambda\, b' \to uncurry\ g\ (a, b'))\ b \\
&= \mathcal{D}\ (\lambda\, a' \to g\ a'\ b)\ a \triangledown \mathcal{D}\ (\lambda\, b' \to g\ a\ b')\ b && \text{-- } uncurry \text{ on functions} \\
&= \mathcal{D}\ (at\ b \circ g)\ a \triangledown \mathcal{D}\ (g\ a)\ b && \text{-- } at \text{ definition and } \eta \text{ reduction} \\
&= \mathcal{D}\ (at\ b)\ (g\ a) \circ \mathcal{D} g\ a \triangledown \mathcal{D}\ (g\ a)\ b && \text{-- chain rule (Theorem 1)} \\
&= at\ b \circ \mathcal{D} g\ a \triangledown \mathcal{D}\ (g\ a)\ b && \text{-- linearity of } at
\end{aligned}
$$

---

[24] The cocartesian law $h = h \circ inl \triangledown h \circ inr$ is dual to the cartesian law $h = exl \circ h \triangle exr \circ h$ (Gibbons, 2002).

## A.3  Corollary 4.3

$$
\begin{aligned}
&\mathcal{D}\ eval\ (f, a) \\
&= \mathcal{D}_l\ eval\ (f, a) \triangledown \mathcal{D}_r\ eval\ (f, a) && \text{-- Lemma 4} \\
&= \mathcal{D}\ (eval \circ (, a))\ f \triangledown \mathcal{D}\ (eval \circ (f, ))\ a && \text{-- } \mathcal{D}_l \text{ and } \mathcal{D}_r \text{ alternative definitions} \\
&= \mathcal{D}\ (at\ a)\ f \triangledown \mathcal{D}\ f\ a && \text{-- } eval \text{ on functions; } at \text{ definition} \\
&= at\ a \triangledown \mathcal{D}\ f\ a && \text{-- linearity of } at\ a \\
&= \lambda\ (df, dx) \rightarrow df\ a + \mathcal{D}\ f\ a\ dx && \text{-- } (\triangledown)\ on\ linear\ maps; \ at \text{ definition}
\end{aligned}
$$

Alternatively, calculate $\mathcal{D}\ eval$ via *uncurry*:

$$
\begin{aligned}
&\mathcal{D}\ eval\ (f, a) \\
&= \mathcal{D}\ (uncurry\ id)\ (f, a) && \text{-- } eval = uncurry\ id \\
&= at\ a \circ \mathcal{D}\ id\ a \triangledown \mathcal{D}\ (id\ f)\ a && \text{-- Corollary 4.2} \\
&= at\ a \circ id \triangledown \mathcal{D}\ f\ a && \text{-- } id \text{ linearity} \\
&= at\ a \triangledown \mathcal{D}\ f\ a && \text{-- } id \text{ as identity}
\end{aligned}
$$

## A.4  Lemma 5

$$
\begin{aligned}
&fork_F\ (\lambda\ b \rightarrow \mathcal{D}\ (at\ b \circ g)\ a) \\
&= \lambda\ da\ b \rightarrow \mathcal{D}\ (at\ b \circ g)\ a\ da && \text{-- } fork_F \text{ definition} \\
&= \lambda\ da\ b \rightarrow (\mathcal{D}\ (at\ b)\ (g\ a) \circ \mathcal{D}\ g\ a)\ da && \text{-- chain rule (Theorem 1)} \\
&= \lambda\ da\ b \rightarrow (at\ b \circ \mathcal{D}\ g\ a)\ da && \text{-- } at\ b \text{ linearity} \\
&= \lambda\ da\ b \rightarrow at\ b\ (\mathcal{D}\ g\ a\ da) && \text{-- } (\circ) \text{ on functions} \\
&= \lambda\ da\ b \rightarrow \mathcal{D}\ g\ a\ da\ b && \text{-- } at \text{ definition} \\
&= \mathcal{D}\ g\ a && \text{-- } \eta \text{ reduction (twice)}
\end{aligned}
$$

## A.5  Corollary 5.1

$$
\begin{aligned}
&\mathcal{D}\ (curry\ f)\ a \\
&= fork_F\ (\lambda\ b \rightarrow \mathcal{D}\ (at\ b \circ curry\ f))\ a && \text{-- Lemma 5} \\
&= fork_F\ (\lambda\ b \rightarrow \mathcal{D}\ (\lambda\ a \rightarrow at\ b\ (curry\ f\ a)))\ a && \text{-- } (\circ) \text{ on functions} \\
&= fork_F\ (\lambda\ b \rightarrow \mathcal{D}\ (\lambda\ a \rightarrow curry\ f\ a\ b))\ a && \text{-- } at \text{ definition} \\
&= fork_F\ (\lambda\ b \rightarrow \mathcal{D}\ (\lambda\ a \rightarrow f\ (a, b)))\ a && \text{-- } curry \text{ on functions} \\
&= fork_F\ (\lambda\ b \rightarrow \mathcal{D}\ (f \circ (, b)))\ a && \text{-- } (, b) \text{ definition} \\
&= fork_F\ (\lambda\ b \rightarrow \mathcal{D}_l\ f\ (a, b)) && \text{-- } \mathcal{D}_l \text{ definition} \\
&= fork_F\ (\mathcal{D}_l\ f \circ (a, )) && \text{-- } (a, ) \text{ definition}
\end{aligned}
$$

## A.6  Lemma 9

The functions $wrap_o$ and $wrap_o^{-1}$ form an isomorphism:

$$
\begin{aligned}
&wrap_o^{-1}\ (wrap_o\ f) \\
&= wrap_o^{-1}\ (o \circ f \circ o^{-1}) && \text{-- } wrap_o \text{ definition} \\
&= o^{-1} \circ (o \circ f \circ o^{-1}) \circ o && \text{-- } wrap_o^{-1} \text{ definition} \\
&= (o^{-1} \circ o) \circ f \circ (o^{-1} \circ o) && \text{-- } (\circ) \text{ associativity} \\
&= id \circ f \circ id && \text{-- } o^{-1} \circ o = id \\
&= f && \text{-- } id \text{ is identity for } (\circ)
\end{aligned}
$$

$$
\begin{aligned}
&wrap_o\ (wrap_o^{-1}\ h) \\
&= wrap_o\ (o^{-1} \circ h \circ o) && \text{-- } wrap_o^{-1} \text{ definition}
\end{aligned}
$$

$$
\begin{aligned}
&= o \circ (o^{-1} \circ h \circ o) \circ o^{-1} && \text{-- } wrap_o \text{ definition}\\
&= (o \circ o^{-1}) \circ h \circ (o \circ o^{-1}) && \text{-- } (\circ) \text{ associativity}\\
&= id \circ h \circ id && \text{-- } o \circ o^{-1} = id\\
&= h && \text{-- } id \text{ is identity for } (\circ)
\end{aligned}
$$

Linearity of $wrap_o$ and $wrap_o^{-1}$ follows from two facts:

- $(\circ f)$ is linear for all $f$.

- $(g \circ)$ is linear for all *linear* $g$.

Proof: exercise.

## A.7   Lemma 10

The functions $\mathcal{D}_o$ and $\mathcal{D}_o^{-1}$ form an isomorphism:

$$
\begin{aligned}
&\mathcal{D}_o^{-1} \circ \mathcal{D}_o\\
&= wrap_o^{-1} \circ \hat{\mathcal{D}}^{-1} \circ \hat{\mathcal{D}} \circ wrap_o && \text{-- } \mathcal{D}_o^{-1} \text{ and } \mathcal{D}_o \text{ definitions}\\
&= wrap_o^{-1} \circ wrap_o && \text{-- } \hat{\mathcal{D}}^{-1} \circ \hat{\mathcal{D}} = id\\
&= id && \text{-- } wrap_o^{-1} \circ wrap_o = id
\end{aligned}
$$

$$
\begin{aligned}
&\mathcal{D}_o \circ \mathcal{D}_o^{-1}\\
&= \hat{\mathcal{D}} \circ wrap_o \circ wrap_o^{-1} \circ \mathcal{D}_o^{-1} && \text{-- } \mathcal{D}_o \text{ and } \mathcal{D}_o^{-1} \text{ definitions}\\
&= \hat{\mathcal{D}} \circ \mathcal{D}_o^{-1} && \text{-- } wrap_o^{-1} \circ wrap_o = id\\
&= id && \text{-- } \hat{\mathcal{D}}^{-1} \circ \hat{\mathcal{D}} = id
\end{aligned}
$$

Linearity of $\mathcal{D}_o$ and $\mathcal{D}_o^{-1}$ follows from linearity of $\hat{\mathcal{D}}$ and $\hat{\mathcal{D}}^{-1}$ and Lemma 9.

## A.8   Lemma 11

The proof that $wrap_o$ is a cartesian functor mainly exploit the regular structure of $o$ and $o^{-1}$:

$$
\begin{aligned}
&wrap_o\ id\\
&= o \circ id \circ o^{-1} && \text{-- } wrap_o \text{ definition}\\
&= o \circ o^{-1} && \text{-- } id \text{ is identity for } (\circ)\\
&= id && \text{-- } o \circ o^{-1} = id
\end{aligned}
$$

$$
\begin{aligned}
&wrap_o\ (g \circ f)\\
&= o \circ g \circ f \circ o^{-1} && \text{-- } wrap_o \text{ definition}\\
&= o \circ g \circ o^{-1} \circ o \circ f \circ o^{-1} && \text{-- } o^{-1} \circ o = id\\
&= (o \circ g \circ o^{-1}) \circ (o \circ f \circ o^{-1}) && \text{-- } \hat{\mathcal{D}} \text{ is a functor}\\
&= wrap_o\ g \circ wrap_o\ f && \text{-- } wrap_o \text{ definition}
\end{aligned}
$$

$$
\begin{aligned}
&wrap_o\ (f \times g) && \text{-- } \mathcal{D}_o \text{ definition}\\
&= o \circ (f \times g) \circ o^{-1} && \text{-- } wrap_o \text{ definition}\\
&= (o \times o) \circ (f \times g) \circ (o^{-1} \times o^{-1}) && \text{-- } o \text{ on products}\\
&= o \circ f \circ o^{-1} \times o \circ g \circ o^{-1} && \text{-- monoidal category law}\\
&= o \circ f \circ o^{-1} \times o \circ g \circ o^{-1} && \text{-- } \hat{\mathcal{D}} \text{ is a monoidal functor}\\
&= wrap_o\ f \times wrap_o\ g && \text{-- } wrap_o \text{ definition}
\end{aligned}
$$

$$
\begin{aligned}
&wrap_o\ exl && \text{-- } \mathcal{D}_o \text{ definition}\\
&= o \circ exl \circ o^{-1} && \text{-- } wrap_o \text{ definition}\\
&= o \circ exl \circ (o^{-1} \times o^{-1}) && \text{-- } o^{-1} \text{ on products}\\
&= o \circ o^{-1} \circ exl && \text{-- } exl \circ (f \times g) = f \circ exl \text{ for cartesian categories}\\
&= exl && \text{-- } o \circ o^{-1} = id
\end{aligned}
$$

$$wrap_o \; exr$$
$$= o \circ exr \circ o^{-1} \qquad\qquad\qquad \text{-- } wrap_o \text{ definition}$$
$$= o \circ exr \circ (o^{-1} \times o^{-1}) \qquad \text{-- } o^{-1} \text{ on products}$$
$$= o \circ o^{-1} \circ exr \qquad\qquad \text{-- } exr \circ (f \times g) = g \circ exr \text{ for cartesian categories}$$
$$= exr \qquad\qquad\qquad\qquad \text{-- } o \circ o^{-1} = id$$

(first line: $wrap_o \; exr$ -- $\mathcal{D}_o$ definition)

$$wrap_o \; dup$$
$$= o \circ dup \circ o^{-1} \qquad\qquad\qquad \text{-- } wrap_o \text{ definition}$$
$$= o \circ (o^{-1} \times o^{-1}) \circ dup \qquad \text{-- } dup \circ f = (f \times f) \circ dup \text{ for cartesian categories}$$
$$= o \circ o^{-1} \circ dup \qquad\qquad \text{-- } o^{-1} \text{ on products}$$
$$= dup \qquad\qquad\qquad\qquad \text{-- } o \circ o^{-1} = id$$

(first line: $wrap_o \; dup$ -- $\mathcal{D}_o$ definition)

## A.9   Lemma 13

$$wrap_o \; (curry \; f)$$
$$= o \circ curry \; f \circ o^{-1} \qquad\qquad \text{-- } wrap_o \text{ definition}$$
$$= \mathcal{D}_o \circ curry \; f \circ o^{-1} \qquad\quad \text{-- } o \text{ on functions}$$
$$= \hat{\mathcal{D}} \circ wrap_o \circ curry \; f \circ o^{-1} \quad \text{-- } \mathcal{D}_o \text{ definition}$$
$$= \hat{\mathcal{D}} \circ curry \; (wrap_o \; f) \qquad\quad \text{-- below}$$

For this last step,

$$wrap_o \circ curry \; f \circ o^{-1}$$
$$= \lambda \; a \to wrap_o \; (curry \; f \; (o^{-1} \; a)) \qquad\qquad \text{-- } \eta \text{ expansion}$$
$$= \lambda \; a \to o \circ curry \; f \; (o^{-1} \; a) \circ o^{-1} \qquad\quad \text{-- } wrap_o \text{ definition}$$
$$= \lambda \; a \; b \to o \; (curry \; f \; (o^{-1} \; a) \; (o^{-1} \; b)) \quad \text{-- } \eta \text{ expansion}$$
$$= \lambda \; a \; b \to o \; (f \; (o^{-1} \; a, \; o^{-1} \; b)) \qquad\quad \text{-- } curry \text{ on functions}$$
$$= \lambda \; a \; b \to o \; (f \; (o^{-1} \; (a, b))) \qquad\qquad \text{-- } o^{-1} \text{ on pairs}$$
$$= \lambda \; a \; b \to wrap_o \; f \; (a, b) \qquad\qquad\quad \text{-- } wrap_o \text{ definition}$$
$$= curry \; (wrap_o \; f) \qquad\qquad\qquad\qquad \text{-- } curry \text{ on functions}$$

Equivalently, $curry \; (wrap_o \; f) = \hat{\mathcal{D}}^{-1} \circ wrap_o \; (curry \; f)$.[25]

## A.10   $\mathcal{D}_o$ and $eval$

The homomorphism equation is $eval = \mathcal{D}_o \; eval$. Simplifying the RHS,

$$\mathcal{D}_o \; eval$$
$$= \hat{\mathcal{D}} \; (wrap_o \; eval) \qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \text{-- } \mathcal{D}_o \text{ definition}$$
$$= \hat{\mathcal{D}} \; (o \circ eval \circ o^{-1}) \qquad\qquad\qquad\qquad\qquad\qquad\quad \text{-- } wrap_o \text{ definition}$$
$$= \hat{\mathcal{D}} \; (o \circ eval \circ (\mathcal{D}_o^{-1} \times o^{-1})) \qquad\qquad\qquad\qquad \text{-- } o^{-1} \text{ on } (a \to b) \times a$$
$$= \hat{\mathcal{D}} \; (\lambda \, (\hat{f}, a) \to (o \circ eval \circ (\mathcal{D}_o^{-1} \times o^{-1})) \, (\hat{f}, a)) \qquad \text{-- } \eta \text{ expansion}$$
$$= \hat{\mathcal{D}} \; (\lambda \, (\hat{f}, a) \to o \; (eval \; (\mathcal{D}_o^{-1} \; \hat{f}, \; o^{-1} \; a))) \qquad\quad \text{-- } (\circ) \text{ and } (\times) \text{ on functions}$$
$$= \hat{\mathcal{D}} \; (\lambda \, (\hat{f}, a) \to o \; (\mathcal{D}_o^{-1} \; \hat{f} \; (o^{-1} \; a))) \qquad\qquad\quad \text{-- } eval \text{ on functions}$$
$$= \hat{\mathcal{D}} \; (\lambda \, (\hat{f}, a) \to o \; (wrap_o^{-1} \; (\hat{\mathcal{D}}^{-1} \; \hat{f}) \; (o^{-1} \; a))) \qquad \text{-- } \mathcal{D}_o^{-1} \text{ definition}$$
$$= \hat{\mathcal{D}} \; (\lambda \, (\hat{f}, a) \to o \; ((o^{-1} \circ \hat{\mathcal{D}}^{-1} \; \hat{f} \circ o) \; (o^{-1} \; a))) \quad \text{-- } wrap_o^{-1} \text{ definition}$$
$$= \hat{\mathcal{D}} \; (\lambda \, (\hat{f}, a) \to o \; (o^{-1} \; (\hat{\mathcal{D}}^{-1} \; \hat{f} \; (o \; (o^{-1} \; a))))) \qquad \text{-- } (\circ) \text{ on functions}$$
$$= \hat{\mathcal{D}} \; (\lambda \, (\hat{f}, a) \to \hat{\mathcal{D}}^{-1} \; \hat{f} \; a) \qquad\qquad\qquad\qquad\quad \text{-- } o \circ o^{-1} = id$$
$$= \hat{\mathcal{D}} \; (uncurry \; \hat{\mathcal{D}}^{-1}) \qquad\qquad\qquad\qquad\qquad\qquad \text{-- } uncurry \text{ on functions}$$
$$= D \; (\lambda \, (\hat{f}, a) \to (uncurry \; \hat{\mathcal{D}}^{-1} \; (\hat{f}, a), \mathcal{D} \; (uncurry \; \hat{\mathcal{D}}^{-1}) \; (\hat{f}, a))) \quad \text{-- } \hat{\mathcal{D}} \text{ definition}$$

---

[25][Maybe this form will help simplify another proof.]

$$= D \ (\lambda \ (\hat{f}, a) \to (\hat{\mathcal{D}}^{-1} \ \hat{f} \ a, \mathcal{D} \ (uncurry \ \hat{\mathcal{D}}^{-1}) \ (\hat{f}, a)))$$    -- *uncurry* on functions

$$= D \ (\lambda \ (\hat{f}, a) \to (\hat{\mathcal{D}}^{-1} \ \hat{f} \ a, at \ a \circ \mathcal{D} \ \hat{\mathcal{D}}^{-1} \ \hat{f} \ \triangledown \mathcal{D} \ (\hat{\mathcal{D}}^{-1} \ \hat{f}) \ a))$$

$$= D \ (\lambda \ (\hat{f}, a) \to (\hat{\mathcal{D}}^{-1} \ \hat{f} \ a, at \ a \circ \hat{\mathcal{D}}^{-1} \ \triangledown \mathcal{D} \ (\hat{\mathcal{D}}^{-1} \ \hat{f}) \ a))$$    -- $\hat{\mathcal{D}}^{-1}$ linearity

Now note that

$$\hat{f}$$
$$= \hat{\mathcal{D}} \ (\hat{\mathcal{D}}^{-1} \ \hat{f}) \qquad\qquad \text{-- } \hat{\mathcal{D}} \circ \hat{\mathcal{D}}^{-1} = id$$
$$= D \ (\hat{\mathcal{D}}^{-1} \ \hat{f} \triangle \mathcal{D} \ (\hat{\mathcal{D}}^{-1} \ \hat{f})) \quad \text{-- } \hat{\mathcal{D}} \text{ definition}$$

Letting $D \ h = \hat{f}$, we have

$$h \ a = (\hat{\mathcal{D}}^{-1} \ \hat{f} \triangle \mathcal{D} \ (\hat{\mathcal{D}}^{-1} \ \hat{f})) \ a$$
$$= (\hat{\mathcal{D}}^{-1} \ \hat{f} \ a, \mathcal{D} \ (\hat{\mathcal{D}}^{-1} \ \hat{f}) \ a)$$

A bit of refactoring then replaces $\hat{\mathcal{D}}^{-1} \ \hat{f} \ a$ and (the noncomputable) $\mathcal{D} \ (\hat{\mathcal{D}}^{-1} \ \hat{f} \ a)$, yielding a *computable* form:

$$\mathcal{D}_o \ eval$$
$$= \dots$$
$$= D \ (\lambda \ (D \ h, a) \to \textbf{let} \ (b, f') = h \ a \ \textbf{in} \ (b, at \ a \circ \hat{\mathcal{D}}^{-1} \ \triangledown f'))$$

Since this calculation was fairly involved, let's get a sanity check on the types in the final form:

$$(D \ h, a) : O \ ((a \to b) \times a)$$
$$: D \ (O \ a) \ (O \ b) \times O \ a$$
$$D \ h \quad : D \ (O \ a) \ (O \ b)$$
$$a \ : O \ a$$
$$h \quad : O \ a \to O \ b \times (O \ a \multimap O \ b)$$
$$h \ a \quad : O \ b \times (O \ a \multimap O \ b)$$

$$(b, f') \quad : O \ b \times (O \ a \multimap O \ b)$$
$$b \qquad : O \ b$$
$$f' \qquad : O \ a \multimap O \ b$$

$$\hat{\mathcal{D}}^{-1} \qquad\qquad\qquad\qquad\qquad : D \ (O \ a) \ (O \ b) \multimap (O \ a \to O \ b)$$
$$at \ a \qquad\qquad\qquad : (O \ a \to O \ b) \multimap O \ b$$
$$at \ a \circ \hat{\mathcal{D}}^{-1} \qquad : D \ (O \ a) \ (O \ b) \multimap O \ b$$
$$at \ a \circ \hat{\mathcal{D}}^{-1} \ \triangledown f' \ : D \ (O \ a) \ (O \ b) \times O \ a \multimap O \ b$$
$$(b, at \ a \circ \hat{\mathcal{D}}^{-1} \ \triangledown f') \ : O \ b \times (D \ (O \ a) \ (O \ b) \times O \ a \multimap O \ b)$$
$$\lambda \ (D \ h, a) \to \dots \textbf{in} \ (b, at \ a \circ \hat{\mathcal{D}}^{-1} \ \triangledown f') \ : O \ ((a \to b) \times a) \to O \ b \times (D \ (O \ a) \ (O \ b) \times O \ a \multimap O \ b)$$
$$D \ (\lambda \ (D \ h, a) \to \dots \textbf{in} \ (b, at \ a \circ \hat{\mathcal{D}}^{-1} \ \triangledown f')) : D \ (O \ ((a \to b) \times a)) \ (O \ b)$$

$$eval : (a \to b) \times a \to b$$
$$\mathcal{D}_o \ eval : D \ (O \ ((a \to b) \times a)) \ (O \ b)$$

## A.11   Lemma 15

Letting $f_o = wrap_o \ f$,

$$\mathcal{D}_o \ (curry \ f)$$
$$= \hat{\mathcal{D}} \ (wrap_o \ (curry \ f)) \qquad\qquad \text{-- } \mathcal{D}_o \text{ definition}$$
$$= \hat{\mathcal{D}} \ (\hat{\mathcal{D}} \circ curry \ (wrap_o \ f)) \qquad \text{-- Lemma 13}$$
$$= \hat{\mathcal{D}} \ (\hat{\mathcal{D}} \circ curry \ f_o) \qquad\qquad \text{-- } f_o \text{ definition}$$

$$
\begin{aligned}
&= D\left((\hat{\mathcal{D}} \circ \mathit{curry}\ f_o) \bigtriangleup \mathcal{D}\left(\hat{\mathcal{D}} \circ \mathit{curry}\ f_o\right)\right) && \text{-- } \hat{\mathcal{D}}\text{ definition} \\
&= D\left(\lambda\, a \to \hat{\mathcal{D}}\left(\mathit{curry}\ f_o\ a\right), \mathcal{D}\left(\hat{\mathcal{D}} \circ \mathit{curry}\ f_o\right) a\right) && \text{-- } (\bigtriangleup)\text{ definition} \\
&= D\left(\lambda\, a \to \hat{\mathcal{D}}\left(\mathit{curry}\ f_o\ a\right), \hat{\mathcal{D}} \circ \mathcal{D}\left(\mathit{curry}\ f_o\right) a\right) && \text{-- chain rule; linearity of } \hat{\mathcal{D}} \\
&= D\left(\lambda\, a \to \hat{\mathcal{D}}\left(\mathit{curry}\ f_o\ a\right), \hat{\mathcal{D}} \circ \mathit{fork}_F\left(\mathcal{D}_l\ f_o \circ (a,)\right)\right) && \text{-- Theorem 2}
\end{aligned}
$$

Now, separately simplify the two main parts of this last form.

$$
\begin{aligned}
&\hat{\mathcal{D}}\left(\mathit{curry}\ f_o\ a\right) \\
&= D\left(\lambda\, b \to \left(f_o\ (a, b), \mathcal{D}_r\ f_o\ (a, b)\right)\right) && \text{-- } \hat{\mathcal{D}}\text{ definition and Lemma 4}
\end{aligned}
$$

$$
\begin{aligned}
&\hat{\mathcal{D}} \circ \mathit{fork}_F\left(\mathcal{D}_l\ f_o \circ (a,)\right) \\
&= \hat{\mathcal{D}} \circ \left(\lambda\, da\ b \to \left(\mathcal{D}_l\ f_o \circ (a,)\right) b\ da\right) && \text{-- } \mathit{fork}_F\text{ definition} \\
&= \hat{\mathcal{D}} \circ \left(\lambda\, da\ b \to \mathcal{D}_l\ f_o\ (a, b)\ da\right) && \text{-- } (\circ)\text{ on functions} \\
&= \lambda\, da \to \hat{\mathcal{D}}\left(\lambda\, b \to \mathcal{D}_l\ f_o\ (a, b)\ da\right) && \text{-- } (\circ)\text{ on functions} \\
&= \lambda\, da \to \hat{\mathcal{D}}\left(\lambda\, b \to at\ da\ \left(\mathcal{D}_l\ f_o\ (a, b)\right)\right) && \text{-- } at\text{ definition} \\
&= \lambda\, da \to \hat{\mathcal{D}}\left(at\ da \circ \mathcal{D}_l\ f_o \circ (a,)\right) && \text{-- } (\circ)\text{ on functions} \\
&= \lambda\, da \to D\left(\lambda\, b \to \left(\left(at\ da \circ \mathcal{D}_l\ f_o \circ (a,)\right) b, \mathcal{D}\left(at\ da \circ \mathcal{D}_l\ f_o \circ (a,)\right) b\right)\right) && \text{-- } \hat{\mathcal{D}}\text{ definition} \\
&= \lambda\, da \to D\left(\lambda\, b \to \left(\mathcal{D}_l\ f_o\ (a, b)\ da, \mathcal{D}\left(at\ da \circ \mathcal{D}_l\ f_o \circ (a,)\right) b\right)\right) && \text{-- } (\circ)\text{ on functions}
\end{aligned}
$$

Now simplify the remaining differentiated composition:

$$
\begin{aligned}
&\mathcal{D}\left(at\ da \circ \mathcal{D}_l\ f_o \circ (a,)\right) b \\
&= at\ da \circ \mathcal{D}\left(\mathcal{D}_l\ f_o \circ (a,)\right) b && \text{-- chain rule; linearity of } at\ da \\
&= at\ da \circ \mathcal{D}_r\left(\mathcal{D}_l\ f_o\right) (a, b) && \text{-- Lemma 4}
\end{aligned}
$$

Putting the pieces back together,

$$
\begin{aligned}
&\mathcal{D}_o\left(\mathit{curry}\ f\right) = \\
&\quad D\left(\lambda\, a \to \left(D\left(\lambda\, b \to \left(f_o\ (a, b), \mathcal{D}_r\ f_o\ (a, b)\right)\right)\right.\right. \\
&\quad\quad\left.\left. , \lambda\, da \to D\left(\lambda\, b \to \left(\mathcal{D}_l\ f_o\ (a, b)\ da, at\ da \circ \mathcal{D}_r\left(\mathcal{D}_l\ f_o\right) (a, b)\right)\right)\right)\right)
\end{aligned}
$$

## A.12 Lemma 16

To show that $comp = uncurry\ (\circ)$ is bilinear, we can show that it is linear in each argument, which is to say $curry\ comp\ g = (g\ \circ)$ and $curry'\ comp\ f = (\circ\ f)$ are linear for all $g$ and $f$.

First, $(\circ\ f)$ is linear for *any* function $f$ (not just linear):

$$
\begin{aligned}
&(\circ\ f)\left(g + g'\right) \\
&= \left(g + g'\right) \circ f && \text{-- left section definition} \\
&= \lambda\, a \to \left(g + g'\right) (f\ a) && \text{-- } \eta\text{ expansion} \\
&= \lambda\, a \to g\ (f\ a) + g'\ (f\ a) && \text{-- addition on functions} \\
&= \left(\lambda\, a \to g\ (f\ a)\right) + \left(\lambda\, a \to g'\ (f\ a)\right) && \text{-- addition on functions} \\
&= \left(g \circ f\right) + \left(g' \circ f\right) && \text{-- } (\circ)\text{ on functions} \\
&= (\circ\ f)\ g + (\circ\ f)\ g' && \text{-- left section definition}
\end{aligned}
$$

$$
\begin{aligned}
&(\circ\ f)\ (s \cdot g) \\
&= (s \cdot g) \circ f && \text{-- left section definition} \\
&= \lambda\, a \to (s \cdot g)\ (f\ a) && \text{-- } (\circ)\text{ on functions} \\
&= \lambda\, a \to s \cdot g\ (f\ a) && \text{-- scaling on functions} \\
&= s \cdot \left(\lambda\, a \to g\ (f\ a)\right) && \text{-- scaling on functions} \\
&= s \cdot (g \circ f) && \text{-- } (\circ)\text{ on functions} \\
&= s \cdot (\circ\ f)\ g && \text{-- left section definition}
\end{aligned}
$$

Second, $(g \circ)$ is linear for any *linear* functions $g$:

$$
\begin{aligned}
&(g \circ) \; (f + f') \\
&= g \circ (f + f') && \text{-- right section definition} \\
&= \lambda\, a \to g\; ((f + f')\; a) && \text{-- } \eta \text{ expansion} \\
&= \lambda\, a \to g\; (f\; a + f'\; a) && \text{-- addition on functions} \\
&= \lambda\, a \to g\; (f\; a) + g\; (f'\; a) && \text{-- linearity of } g \\
&= (\lambda\, a \to g\; (f\; a)) + (\lambda\, a \to g\; (f'\; a)) && \text{-- addition on functions} \\
&= (g \circ f) + (g \circ f') && \text{-- } (\circ) \text{ on functions} \\
&= (g \circ)\; f + (g \circ)\; f' && \text{-- right section definition}
\end{aligned}
$$

$$
\begin{aligned}
&(g \circ) \; (s \cdot f) \\
&= g \circ (s \cdot f) && \text{-- right section definition} \\
&= \lambda\, a \to g\; ((s \cdot f)\; a) && \text{-- } (\circ) \text{ on functions} \\
&= \lambda\, a \to g\; (s \cdot f\; a) && \text{-- scaling on functions} \\
&= \lambda\, a \to s \cdot g\; (f\; a) && \text{-- linearity of } g \\
&= s \cdot (g \circ f) && \text{-- scaling on functions} \\
&= s \cdot (g \circ)\; f && \text{-- right section definition}
\end{aligned}
$$

## A.13  Lemma 17

Given any bilinear function $h$,

a. *curry h a* is linear for all linear functions $g$:

$$
\begin{aligned}
&curry\; h\; a\; (b + b') \\
&= h\; (a, b + b') && \text{-- } curry \text{ on functions} \\
&= h\; (a, b) + h\; (a, b') && \text{-- bilinearity of } h \\
&= curry\; h\; a\; b + curry\; h\; a\; b' && \text{-- } curry \text{ on functions}
\end{aligned}
$$

$$
\begin{aligned}
&curry\; h\; a\; (s \cdot b) \\
&= h\; (a, s \cdot b) && \text{-- } curry \text{ on functions} \\
&= s \cdot h\; (a, b) && \text{-- bilinearity of } h
\end{aligned}
$$

b. *curry' h b* is linear for all functions $b$: Proof similar to *curry h a*.

c. *curry h* and *curry' h* are linear:

$$
\begin{aligned}
&curry\; h\; (a + a') \\
&= \lambda\, b \to curry\; h\; (a + a')\; b && \text{-- } \eta \text{ expansion} \\
&= \lambda\, b \to h\; (a + a', b) && \text{-- } curry \text{ on functions} \\
&= \lambda\, b \to h\; (a, b) + h\; (a', b) && \text{-- bilinearity of } h \\
&= (\lambda\, b \to h\; (a, b)) + (\lambda\, b \to h\; (a', b)) && \text{-- addition on functions} \\
&= curry\; h\; a + curry\; h\; a' && \text{-- } curry \text{ on functions}
\end{aligned}
$$

$$
\begin{aligned}
&curry\; h\; (s \cdot a) \\
&= \lambda\, b \to curry\; h\; (s \cdot a)\; b && \text{-- } \eta \text{ expansion} \\
&= \lambda\, b \to h\; (s \cdot a', b) && \text{-- } curry \text{ on functions} \\
&= \lambda\, b \to s \cdot h\; (a, b) && \text{-- bilinearity of } h \\
&= s \cdot (\lambda\, b \to h\; (a, b)) && \text{-- scaling on functions} \\
&= s \cdot curry\; h\; a && \text{-- } curry \text{ on functions}
\end{aligned}
$$

Similarly for *curry' h*.

d. $\mathcal{D}\ h$ is linear:

$$
\begin{array}{lll}
& \mathcal{D}\ h\ ((a, b) + (a', b')) & \\
= & \mathcal{D}\ h\ (a + a', b + b') & \text{-- (+) on functions} \\
= & h \circ (, b + b') \triangledown h \circ (a + a') & \text{-- Corollary 4.1} \\
= & \lambda\ (da, db) \rightarrow h\ (da, b + b') + h\ (a + a', db) & \text{-- } (\triangledown) \text{ on functions} \\
= & \lambda\ (da, db) \rightarrow h\ (da, b) + h\ (da, b') + h\ (a, db) + h\ (a', db) & \text{-- bilinearity of } h \\
= & \lambda\ (da, db) \rightarrow h\ (da, b) + h\ (a, db) + h\ (da, b') + h\ (a', db) & \text{-- commutativity of (+)} \\
= & (\lambda\ (da, db) \rightarrow h\ (da, b\ ) + h\ (a\ , db)) + & \\
& (\lambda\ (da, db) \rightarrow h\ (da, b') + h\ (a', db)) & \text{-- (+) on functions} \\
= & \mathcal{D}\ h\ (a, b) + \mathcal{D}\ h\ (a', b') & \text{-- Corollary 4.1}
\end{array}
$$

Similarly for scaling.

## A.14  Lemma 18

a. Sequential composition:

$$
\begin{array}{lll}
& \mathcal{D}\ (g \circ f) & \\
= & \lambda\ a \rightarrow \mathcal{D}\ (g \circ f)\ a & \text{-- } \eta \text{ expansion} \\
= & \lambda\ a \rightarrow \mathcal{D}\ g\ (f\ a) \circ \mathcal{D}\ f\ a & \text{-- chain rule (Theorem 1)} \\
= & \lambda\ a \rightarrow (\mathcal{D}\ g \circ f)\ a \circ \mathcal{D}\ f\ a & \text{-- } (\circ) \text{ on functions} \\
= & \lambda\ a \rightarrow (\circ)\ ((\mathcal{D}\ g \circ f)\ a)\ (\mathcal{D}\ f\ a) & \text{-- alternative notation} \\
= & \lambda\ a \rightarrow uncurry\ (\circ)\ ((\mathcal{D}\ g \circ f)\ a, \mathcal{D}\ f\ a) & \text{-- } uncurry \text{ on functions} \\
= & \lambda\ a \rightarrow comp\ ((\mathcal{D}\ g \circ f)\ a, \mathcal{D}\ f\ a) & \text{-- } comp \text{ definition} \\
= & comp \circ (\lambda\ a \rightarrow ((\mathcal{D}\ g \circ f)\ a, \mathcal{D}\ f\ a)) & \text{-- } (\circ) \text{ on functions} \\
= & comp \circ (\mathcal{D}\ g \circ f \triangle \mathcal{D}\ f) & \text{-- } (\triangle) \text{ definition}
\end{array}
$$

b. Cross:

$$
\begin{array}{lll}
& \mathcal{D}\ (f \times g) & \\
= & \lambda\ (a, b) \rightarrow \mathcal{D}\ (f \times g)\ (a, b) & \text{-- } \eta \text{ expansion} \\
= & \lambda\ (a, b) \rightarrow \mathcal{D}\ f\ a \times \mathcal{D}\ g\ b & \text{-- cross rule (Theorem 2)} \\
= & \lambda\ (a, b) \rightarrow uncurry\ (\times)\ (\mathcal{D}\ f\ a, \mathcal{D}\ g\ b) & \text{-- } uncurry \text{ on functions} \\
= & \lambda\ (a, b) \rightarrow cross\ (\mathcal{D}\ f\ a, \mathcal{D}\ g\ b) & \text{-- } cross \text{ definition} \\
= & \lambda\ (a, b) \rightarrow cross\ ((\mathcal{D}\ f \times \mathcal{D}\ g)\ (a, b)) & \text{-- } (\times) \text{ on functions} \\
= & cross \circ (\mathcal{D}\ f \times \mathcal{D}\ g) & \text{-- } (\circ) \text{ on functions}
\end{array}
$$

c. Fork:

$$
\begin{array}{lll}
& \mathcal{D}\ (f \triangle g) & \\
= & \mathcal{D}\ ((f \times g) \circ dup) & \text{-- cartesian law} \\
= & \lambda\ a \rightarrow \mathcal{D}\ ((f \times g) \circ dup)\ a & \text{-- } \eta \text{ expansion} \\
= & \lambda\ a \rightarrow \mathcal{D}\ (f \times g)\ (dup\ a) \circ \mathcal{D}\ dup\ a & \text{-- chain rule (Theorem 1)} \\
= & \lambda\ a \rightarrow \mathcal{D}\ (f \times g)\ (a, a) \circ \mathcal{D}\ dup\ a & \text{-- } dup \text{ for functions} \\
= & \lambda\ a \rightarrow \mathcal{D}\ f\ a \times \mathcal{D}\ g\ a \circ \mathcal{D}\ dup\ a & \text{-- cross rule (Theorem 2)} \\
= & \lambda\ a \rightarrow \mathcal{D}\ f\ a \times \mathcal{D}\ g\ a \circ dup & \text{-- } dup \text{ linearity} \\
= & \lambda\ a \rightarrow \mathcal{D}\ f\ a \triangle \mathcal{D}\ g\ a & \text{-- cartesian law} \\
= & fork \circ (\mathcal{D}\ f \triangle \mathcal{D}\ g) & \text{-- } fork \text{ definition}
\end{array}
$$

d. A linear function $f$,

$$
\begin{array}{lll}
& \mathcal{D}\ f & \\
= & \lambda\ a \rightarrow \mathcal{D}\ f\ a & \text{-- } \eta \text{ expansion} \\
= & \lambda\ a \rightarrow f & \text{-- Theorem 3} \\
= & const\ f & \text{-- } const \text{ definition}
\end{array}
$$

e. Any function $f : a \times b \to c$,

$$
\begin{aligned}
& \mathcal{D} f \\
&= \lambda\,(a, b) \to \mathcal{D} f\,(a, b) && \text{-- } \eta \text{ expansion} \\
&= \lambda\,(a, b) \to \mathcal{D}_l f\,(a, b) \,\triangledown\, \mathcal{D}_r f\,(a, b) && \text{-- Lemma 4} \\
&= join \circ (\mathcal{D}_l f \,\triangle\, \mathcal{D}_r f) && \text{-- } join \text{ definition}
\end{aligned}
$$

f. A *bilinear* function $f : a \times b \to c$,

$$
\begin{aligned}
& \mathcal{D} f \\
&= \lambda\,(a, b) \to \mathcal{D} f\,(a, b) && \text{-- } \eta \text{ expansion} \\
&= \lambda\,(a, b) \to f \circ (, b) \,\triangledown\, f \circ (a, ) && \text{-- Corollary 4.1} \\
&= \lambda\,(a, b) \to curry'\, f\ b \,\triangledown\, curry\ f\ a && \text{-- section definitions} \\
&= \lambda\,(a, b) \to join\,(curry'\, f\ b,\, curry\ f\ a) && \text{-- } join = uncurry\ (\triangledown) \\
&= \lambda\,(a, b) \to join\,((curry'\, f \times curry\ f)\,(b, a)) && \text{-- } (\times) \text{ on functions} \\
&= \lambda\,(a, b) \to join\,((curry'\, f \times curry\ f)\,(swap\,(a, b))) && \text{-- } swap \text{ on functions} \\
&= join \circ (curry'\, f \times curry\ f) \circ swap && \text{-- } (\circ) \text{ on functions}
\end{aligned}
$$

g. Uncurried composition on linear maps,

$$
\begin{aligned}
& \mathcal{D}\ comp \\
&= join \circ (curry'\ comp \times curry\ comp) \circ swap && \text{-- previous (} comp \text{ is bilinear)} \\
&= join \circ (flip\,(\circ) \times (\circ)) \circ swap && \text{-- } comp \text{ definition}
\end{aligned}
$$