

$\partial P.jl$ terminology and notation guide

Keno Fischer

FOREWORD

The field of automatic differentiation (AD) has a long history, both in the academic literature and in programming practice. However, despite this long history, or perhaps because of it, the terminology for various concepts and operations is neither standard nor consistent. This inconsistency imposes upon new implementers of AD tools a burden to clearly state the meaning of any terminology they might venture to employ. This document is an attempt to do such that for my present work on a new AD tool for the Julia language. It is intended 1) to give contributors and users a common vocabulary to discuss the manifold abstractions required to create a modern AD tools, 2) to explain why certain terminology may have been chosen and 3) to serve implementers as a guide to the proper functioning of the tool when lost deep in the web of details. By necessity, the terminology in this guide is borrowed from several different fields of study (most relevantly category theory and differential geometry). Nevertheless, this guide does not intend to be an introduction to these fields, nor does it assume that readers have had such an introduction. Instead, we try to give some examples and intuition for why the abstractions chosen are of interest to the implementation of automatic differentiation. Lastly, despite using terminology that might be amenable to formalization, this document is not an attempt to do so and as such certain liberties may be taken that would require careful treatment in a more formal treatise. With that out of the way, let's begin:

I. OPTICAL CONSTRUCTIONS

To begin our quest, we shall first put aside the topic of differentiation entirely and study a relatively recent construction from category theory known as an *optic*. The power of the optic construction lies in its ability to combine both a covariant and a contravariant transformation in one abstraction while retaining composability. Following Riley [1], we briefly consider the following explicit description of an optic (those not familiar with the language of category theory are encouraged to skip straight to the diagrams).

Definition 0.1: For some symmetric monoidal category \mathcal{C} , we define the category $\mathbf{Optic}_{\mathcal{C}}$. Whose objects are pairs (A, A') of objects in \mathcal{C} and whose morphisms have representatives given by $\langle l|r \rangle : (A, A') \rightarrow (B, B')$ are pairs (l, r) where $l : A \rightarrow M \otimes B$, $r : M \otimes B' \rightarrow A'$ (the choice of the object M depending on the representative - see Riley for details).

This definition makes manifest the combination of co- and contravariant data. For a representative $\langle l|r \rangle$, l varies covariantly while r varies contravariantly. We additionally

have a “memory” or “residual” object M . This object is not uniquely determined and in fact we shall make good use of that fact in our actual use case later. For now, We briefly exhibit the composition rules for these morphisms. Suppose $\langle l_1|r_1 \rangle : (A, A') \rightarrow (B, B')$, $\langle l_2|r_2 \rangle : (B, B') \rightarrow (C, C')$. Then the sequential composition $\langle l_{12}|r_{12} \rangle = \langle l_1|r_1 \rangle \circ \langle l_2|r_2 \rangle$ may be written explicitly as

$$l_{12}(a) := (m_1 \otimes m_2, c) \text{ where } (m_1, b) = l_1(a), (m_2, c) = l_2(b)$$

$$r_{12}(m_1 \otimes m_2, c') := r_1(m_1, r_2(m_2, c'))$$

While these definitions are perfectly explicit, it might be easier to understand what is going on by considering optics diagrammatically (diagrams taken from Riley [1]). For $\langle l|r \rangle$ we draw:



We compose these optics by inserting one optic in the hole of another, e.g. for the composition we have above, we might write $\langle l_1|r_1 \rangle \circ \langle l_2|r_2 \rangle$



From the diagrammatic view, it is immediately clear that this composition is indeed another optic and what the composition rule for these optics should be:



The reader is encouraged to verify that this diagrammatic derivation of l_{12} and r_{12} matches the algebraic definition above.

A. Encoding Optics in Julia

One of the original applications of optic constructions is the manipulation of deeply nested data structures, so let's take that example and consider how we might represent it in Julia. Optics for this application are often called *lenses*, though since we are just using them as an example here we shall not concern ourselves with a precise definition. We begin by defining an abstract optic type as well as a concrete representation of morphisms as pairs $\langle l|r \rangle$:

Snippet 1

```

1  abstract type AbstractOptic; end
2
3  struct OpticRepr <: AbstractOptic
4      l::Function
5      r::Function
6  end
7  OpticRepr(o::OpticRepr) = o
8  OpticRepr(o::AbstractOptic) =
9      OpticRepr(a->left(o, a),
10         (m, b')->right(o, m, b'))
11
12  left(o::OpticRepr, a) = o.l(a)
13  right(o::OpticRepr, m, b') = o.r(m, b')
14
15  function §(o1::AbstractOptic, o2::AbstractOptic)
16      OpticRepr(
17          function (a)
18              x = left(o1, a)
19              y = left(o2, x[2])
20              (x[1], y[1]), y[2]
21          end,
22          ((m1, m2), c')->right(o1, m1, right(o2, m2, c'))
23      )
24  end
25

```

Given these definitions, we can now define an optic constructor that allows us to modify an immutable value at a given location:

Snippet 2

```

1  SetIndex(idx) = OpticRepr(
2      obj->(obj, getindex(obj, idx)),
3      (obj, update)->setindex(obj, update, idx)
4  )

```

which we may use like so:

```

julia> tup = (1, (2, 3))
julia> o = SetIndex(2) § SetIndex(1)
julia> m, c = left(o, tup)
julia> c
2
julia> right(o, m, c + 100)
(1, (102, 3))

```

However, as written, this definition has a big problem: It gives us the state 'm' explicitly. The type and value of 'm' depends on the representative of the optic, but nothing in our above definition requires the state 'm' to be matched to the optic representative 'o'. To see the representation problem explicitly, suppose we noticed that the state in this representation is somewhat suboptimal:

```

julia> m
((1, (2, 3)), (2, 3))

```

and instead defined a more optimized representation of the same optic:

Snippet 3

```

1  struct FastSetIndex <: AbstractOptic
2      path::NTuple{N, Int} where N
3  end
4  FastSetIndex(i::Int) = FastSetIndex((i,))
5
6  path_getindex(obj, ::Tuple{}) = obj
7  path_getindex(obj, path::Tuple) =
8      path_getindex(getindex(obj, first(path)), tail(path))
9  left(f::FastSetIndex, obj) =
10      (obj, path_getindex(obj, f.path))
11
12  path_setindex(obj, update, (idx,)::Tuple{Int}) =
13      setindex(obj, update, idx)
14  path_setindex(obj, update, path) =
15      setindex(obj,
16         path_setindex(getindex(obj, first(path)),
17             update, tail(path)),
18         first(path))
19  function right(f::FastSetIndex, obj, update)
20      path_setindex(obj, update, f.path)
21  end
22
23  (a::FastSetIndex § b::FastSetIndex) =
24      FastSetIndex((a.path..., b.path...))

```

Our example from above goes through just the same:

```

julia> tup = (1, (2, 3))
julia> of = FastSetIndex(2) § FastSetIndex(1)
julia> m, c = left(of, tup)
julia> c
2
julia> right(of, m, c + 100)
(1, (102, 3))

```

but our state is much smaller:

```

julia> m
(1, (2, 3))

```

With these two different representatives for the same optic, the representative, dependence of 'm' becomes quite apparent:

```

julia> right(of, m, c + 100)
(1, (102, 3))

```

```

julia> right(o, m, c + 100)
ERROR: MethodError: no method matching
setindex(::Int64, ::Tuple{Int64, Int64}, ::Int64)

```

In this example, this might be quite silly, but in general it presents a problem depending on how the representative is chosen. E.g. an optic constructor might perform a randomized search for the most efficient representative, so two calls to the constructor will not necessarily yield the same representative. We would thus prefer to have an interface to our optics that either does not explicitly expose the state to the user, or if it does, tags the state with a particular choice or representative. Here are two (equivalent) interfaces to optics that have this property:

Snippet 4

```

1  # Block style optics interface
2  function (o::AbstractOptic)(f::Function, a)
3      m, b = left(o, a)
4      b' = f(b)
5      right(o, m, b')
6  end
7
8  # Continuation style optics interface
9  function (o::AbstractOptic)(a)

```

```

10     m, b = left(o, a)
11     b, b' -> right(o, m, b')
12 end

```

With these definitions, our update operation from above becomes

```

julia> of(tup) do c
    c + 100
end
(1, (102, 3))

julia> c, back = of(tup);

julia> back(c+100)
(1, (102, 3))

```

The observant reader may at this point remark that written this way, we are simply describing a standard design pattern, found frequently in julia libraries (indeed so frequently that the *do* syntax was partly introduced to support it). For example, the julia standard library includes the *cd* function, which modulo error handling function looks like this:

```

function cd(f::Function, dir::AbstractString)
    old = pwd()
    cd(dir)
    f()
    cd(old)
end

```

While it's not a precise match for the optics encoding as we have defined it (and of course is side-effectful which we have not discussed at all).

```

OpticRepr(a->(pwd(), (cd(a[1]); tail(a))),
(m, b)->cd(m))

```

Have we thus wasted the better part of two pages exploring and giving a name to a notion that is so natural as to be a standard design pattern? No! The key takeaway from this section should be twofold optics: 1) Optics have a nice composition property that allows you to take two optics and combine them into another optic and 2) Optical constructions are quite natural and you are probably already familiar with them.

II. SSA AS A CATEGORY / OPTICS OVER SSA

In Julia's optimizer, functions are represented by their single static assignment (SSA) form. We shall venture to sketch a categorical description of SSA form and then investigate the structure of Optic over such a category.

A. Duplication / Deletion of Values

In our examples, so far every values was used once and exactly once. In SSA form, we may have an arbitrary number of uses. To support this, we extend our category with the diagonal $\Delta_{SSA} : A \rightarrow A \otimes A$ and deletion ($\vdash : A \rightarrow I$) morphisms. Graphically, we have:

$$\begin{array}{c} \text{A} \\ \swarrow \quad \searrow \\ \text{A} \quad \text{A} \end{array} \quad \text{A} \longrightarrow \bullet \quad (4)$$

B. Forward control flow

To account for control flow, we add a coproduct to our category, which we shall denote by \oplus . Note that this is different from the \otimes operation, we have so far considered, and we are thus creating a **bimonoidal** category, with two separate monoidal structures. Graphically, we will indicate \oplus by boxes around the operations in questions, so for example, for the morphism $(f \oplus g) \otimes h : (A \oplus B) \otimes C \rightarrow (A' \oplus B') \otimes C'$, we will draw:

$$\begin{array}{|c|} \hline A \rightarrow A' \\ \hline B \rightarrow B' \\ \hline \end{array} \quad (5)$$

$$C \rightarrow C'$$

Now, we have an operation $\varphi : A \oplus A \rightarrow A$, that joins control flow edges:

$$\begin{array}{|c|} \hline A \\ \hline \varphi \\ \hline A \end{array} \rightarrow A \quad (6)$$

and we will introduce an additional “conditonal branch” operation $br : A \otimes C \rightarrow A \oplus A$ (where we say C is the condition):

$$\begin{array}{c} A \\ \text{---} \blacksquare \nearrow \\ C \text{---} \nwarrow \end{array} \begin{array}{|c|} \hline A \\ \hline A \end{array} \quad (7)$$

We further impose the following *co-predication relation*:

$$\begin{array}{c} A \\ \text{---} \blacksquare \nearrow \\ C \text{---} \nwarrow \\ B \\ \text{---} \blacksquare \nearrow \\ B \end{array} \begin{array}{|c|} \hline A \\ \hline A \\ \hline B \\ \hline B \end{array} \cong A \otimes B \begin{array}{c} \text{---} \blacksquare \nearrow \\ \text{---} \nwarrow \end{array} \begin{array}{|c|} \hline A \otimes B \\ \hline A \otimes B \end{array} \quad (8)$$

As a notational convenience, we write \oplus_c for the coproduct introduced under some abstract condition c . In that case, we may write the co-predication condition as

$$\leftrightarrow_{\oplus} := (A \oplus_c B) \otimes (C \oplus_c D) \cong (A \otimes C) \oplus_c (B \otimes D)$$

To avoid running into technical difficulties, we further allow commuting arbitrary morphisms through the branch operations, such that we can be sure every $A \oplus B$ is always potentially the direct result of a branch operation, such that we may validly introduce abstract condition tags on any such \oplus operation¹.

¹I realize this is a bit handwavy. While I believe this works fine, I would be interested in improvements to make this construction more formal and potentially avoid the reliance on commuting morphisms through the branch operation so as ot facilitate the future addition of side effects.

C. Lifting to the Optic category

1) *Product Structure*: Having defined the structure for our bimonoidal category $SSA_{\mathcal{C}}$, what can we say about the structure of $\text{Optic}_{SSA_{\mathcal{C}}}$? To be explicit, we're still performing the optic construction with respect to the \otimes structure of our category. Let us first look at the product itself. Do we have a diagonal map in Optic_{SSA} ? Perhaps surprisingly, the answer is no! To see this, let's consider what we must construct:

We wish to find a map $\text{Optic}(\Delta_{SSA}) : (A, A') \rightarrow (A, A') \otimes (A, A')$. A representative for such a morphism would have the form $\langle l, r \rangle$: $l : A \rightarrow M \otimes A \otimes A$, $r : M \otimes A' \otimes A' \rightarrow A'$. Finding l is easy, we can just use Δ_{SSA} (letting $M = I$) from the underlying category. However, there is no good way to construct r . A priori, there is no reason to expect any morphism of the shape $A' \otimes A' \rightarrow A'$ to exist, and even if such morphisms do exist, there is no reason to expect them to have the required uniqueness properties. It is important to note that this does not prevent us from lifting Δ_{SSA} to Optic in our applications of interest, it simply means that \otimes is not a categorical product in $\text{Optic}_{SSA_{\mathcal{C}}}$, so we must make a choice of summation morphism in our lifting functor.

2) *Coproduct Structure*: Given our utter disappointment with the product structure, can we have any hope to lift the co-product structure. Yes, we do! First we construct the co-product itself. For two optics $\langle l_1 | r_1 \rangle : (A, A') \rightarrow (B, B')$ with residual M_1 and $\langle l_2 | r_2 \rangle : (C, C') \rightarrow (D, D')$ with residual M_2 , we construct a new optic $\langle l_{12} | r_{12} \rangle$ where

$$\begin{aligned} l_{12} &= (l_1 \oplus l_2) \circ \hookrightarrow_{\oplus} \\ r_{12} &= \hookrightarrow_{\oplus}^{-1} \circ (r_1 \oplus r_2) \end{aligned} \quad (9)$$

Similarly, we can confirm that this is indeed a coproduct by exhibiting the φ morphism. The trick here is to use the residual to carry the branch history information. Perhaps the simplest representative of $\text{Optic}(\varphi) : (A \oplus_a A, A' \oplus_a A') \rightarrow (A, A')$ would be $\langle l | r \rangle$ where

$$\begin{aligned} l &= \Delta_{SSA} \circ (\text{id}_{A \oplus_a A} \otimes \varphi) \\ r &= d \circ (\pi_2 \oplus \pi_2) \end{aligned} \quad (10)$$

where d is the distributive map. While this may be a little abstract, hopefully it becomes clear as a picture:



Note that another representative of the same optic could simply record the active component (e.g. by using the map $(x \rightarrow 1) \oplus (y \rightarrow 2)$ and reintroducing a branch on the RHS).

However, note that (at least the way we chose to set it up) we do not have a unique lifting of the branch operator itself.

III. OPTIC FUNCTORS

One useful way to construct optics is to do so automatically, from simpler descriptions. More formally, we are interested in functors $F : \mathcal{C} \rightarrow \mathbf{Optic}_{\mathcal{D}}$. The requirement for F to be a functor means that composition in \mathcal{C} matches composition in $\mathbf{Optic}_{\mathcal{D}}$, i.e. $F(f) \circ F(g) = F(f \circ g)$. We shall call such functors (whose codomain is an optic category) *optic functors*.

We have already seen an example of this kind of operation! In our definition of *FastSetIndex* above, we represented the optic, as a path of indices. If we wanted to, we could consider this path of indices description a category (e.g. the category in which objects are Julia datatypes and there is morphisms from a datatype to each of its field types). In this case, the optic functor is essentially the definitions of the *left* and *right* generic, functions. Of course, we could have also left off the subtype and instead defined a constructor of *OpticRepr* with the *left* and *right* definitions inline (such a constructor does exist in our above definition, by virtue of the subtype). In that case, we might want to check the functor law like so:

```
julia> o1 = OpticRepr(FastSetIndex{2}() ;
                    FastSetIndex{1}())
julia> o2 = OpticRepr(FastSetIndex{2}() ;
                    OpticRepr(FastSetIndex{1}()))
julia> o1(c->c+100, tup)
(1, (102, 3))

julia> o2(c->c+100, tup)
(1, (102, 3))
```

However, it is import to emphasize that while optic functors are required to preserve composition, when implemented as above, they are **not** required to produce the same representative. Indeed, in our example:

```
julia> o1.l(tup)[1]
(1, (2, 3))

julia> o2.l(tup)[1]
((1, (2, 3)), (2, 3))
```

the two representatives o_1 and o_2 are not identical, even though they represent the same optic.

A. Optic functors on programs

We now leave behind our trusty *setindex* example and consider optic functors from programs. To keep the discussion somewhat concrete, consider the following example: Alice, Bob, and Clara are department heads at Consolidated Widgets Incorporated. One day the anti-trust regulator decides that Consolidated Widgets Inc is a little too consolidated and trisects the company into AliceCo, BobCo and ClaraCo with the three former department heads promoted to the CEOs of their respective companies. They now have a problem: While they know how to manufacture widgets, they now also need to keep track of how much each individual part costs and assign profits back to the appropriate company. Can optic functors help? Yes!

Suppose the program for widget assembly is simply

```
widget(sweat, love) = C(A(sweat), B(love))
```

i.e. Alice and Bob both manufacture subassemblies from raw inputs (sweat and love) and Clara handles the final

assembly). We wish to construct an optic functor \mathbb{E} that takes in the widget construction function and spits out a function that both constructs the widgets and accumulates prices on the left side of the optic, and then distributes the profits on the right side of the optic after the widgets have been sold. Now, we have four function of interest here *widget*, *A*, *B* and *C*. For each of *A*, *B* and *C*, we need to do some manual work. Alice, Bob, and Clara need to do some manual work to figure out what their costs are. We say that *A*, *B* and *C* are *primitive* with respect to this functor. Alright, let's suppose we have our prices (for simplicity 10€ for Alice, 20€ for Bob and 30€ for Clara). We will also give each of them a flat 30% profit margin. We will use the continuation encoding of the optic to specify the value of \mathbb{E} on each of our primitives:

Snippet 4

```

1  €A(price, sweat) = let cost=10;
2    (price + cost, A(sweat)),
3    profit->(profit - cost * 1.3)
4  end
5  €B(price, love) = let cost=20;
6    (price + cost, B(love)),
7    profit->(profit - cost * 1.3)
8  end

```

C is a little more tricky. For *C*, we have two arguments, so on right side of the optic, we need to produce two outputs. There are several ways to model this. We could consider only allowing function with a single input and a single output, together with primitive operations for packing (into a tuple) and unpacking, or we could simply require the output to always be a tuple whose length matches the original number of arguments. Here we pick the latter encoding and write:

Snippet 5

```

1  €C(aprice, bprice, a, b) = let cost=30;
2    (aprice + bprice + cost, C(a, b)),
3    profit->(take(profit, aprice),
4             take(profit, bprice))
5  end

```

Ok, phew that was a bit of work to figure out all the primitives. However, now that this is done, we do get some reward for our efforts. We get the definition of $\mathbb{E}widget$ for free by applying the functor law. Recall that we essentially defined $widget = (A \otimes B) \circ C$. Thus, applying the functor law, we immediately get that $\mathbb{E}widget = (\mathbb{E}A \otimes \mathbb{E}B) \circ \mathbb{E}C$

By tracing through definitions, we can write out an explicit representative in the continuation encoding:

Snippet 6

```

1  function €widget(sweat_cost, love_cost, sweat, love)
2    aprice, a, a* = €A(sweat_cost, sweat)
3    bprice, b, b* = €B(love_cost, love)
4    price, c, c* = €C(aprice, bprice, a, b)
5    price, c, profit->begin
6      (a†, b†) = c*(profit)
7      love† = b*(b†)
8      sweat† = a*(a†)
9      (sweat†, love†)
10   end
11 end

```

Now, this might seem a bit complicated (and it is), but that is precisely the point. Writing this out explicitly is quite tedious, but there are no choices to be made in how it is

written. Once we have defined the optic functor on the requisite primitives, we can make use of the composition law to generate a representative of the functor applied to any function that is a composition for our primitives.

B. Some notation

As we saw in the previous section, the explicit representatives generated by optic functors can be quite complicated. To avoid getting lost, we should make some notational conventions. To begin with, we shall consider input programs to be represented by their SSA form representation. We shall assume that the reader is familiar with SSA form IR (if not, see the Julia devdocs on SSA form IR). Now, let \mathbb{E} be an optic functor. In the context of applying \mathbb{E} to some function *f*, we shall say that *f* is the *primal function* and similarly we shall call the value computed by *f* the *primal value*. From now on, we will generally be working with the continuation encoding of an optic. By abuse of notation, we will write $\mathbb{E}f$ for the function that is the canonical representative of the optic in continuation encoding (canonical here meaning obtained by straightforward application of the functor laws). We shall call this function the *optic function*. We shall call the value returned from the optic function the *focal bundle*. Since we sometimes want to look at it in isolation, we will call the continuation generated by the optic function, the *pullback function*, written f^* for some primal function *f*. Since f^* depends on the representative, unless otherwise specified, f^* shall refer to the pullback function generated by the canonical optic function.

When looking at SSA form, IR, for some SSA value $\%a = A(\dots)$ in the primal function, we will write $(\%a, \%a^*) = \mathbb{E}A(\dots)$ in the optic function. The intended implication here is that $\%a$ in the optic function is the same value as $\%a$ in the primal function. This will generally be the case, as long as the optic primitives obey this constraint. However, we have not required this invariant on primitives in our definition of the optic functor and there are (rare) cases in which it is useful for these to be different. Where the distinction matters, the context should be explicitly specified. Lastly, we look at the pullback function. By the optic composition law, pullbacks behave contravariantly to the primal function, so for a call $c = C(A, B)$ in the primal function, we will have a call $(\%a^\dagger, \%b^\dagger) = \%c^*(\%c^\dagger)$. Where in general an underline indicates that a value was captured from the optic function.

C. Optic functors from SSA IR

When we discussed SSA form in section II, we identified several structural features of the IR that did not lift uniquely to the optic category. In particular, this means that we have some freedom in choosing how these should lift for our optic functor. In particular, the structural aspects of SSA IR that did not have a unique lift were:

- Copy
- Delete
- Conditional Branch

As an exercise, let us go through these and choose an implementation.

1) *Copy*: In our little widget manufacturing example from the previous section, the objects being propagated in the primal functions were in general intermediate products that could only be used. But what happens if we introduce trivially copyable goods into the equation? E.g. suppose Alice and Bob both gain a trivially copyable input (perhaps music they like to listen to, or maybe they both require a license for JuliaPro), and that Dave produces such an input. For simplicity, we will also drop Alice and Bob's other inputs. We have a new program:

```

Snippet 7
1 function €clonewidget(d)
2   d = D()
3   a = A(d)
4   b = B(d)
5   c = C(a, b)
6 end

```

Let us look at some sensible policy choices² we might make. For example, we could decide that David gets to decide what to charge for each copy of 'd':

```

Snippet 8
1 function €clone(dprice, d)
2   (dprice, dprice), (d, d),
3   (profit1, profit2) -> (profit1 + profit2)
4 end

```

or maybe everybody pays a fixed fee, independent of David's costs:

```

Snippet 9
1 function €clone(dprice, d)
2   (1, 1), (d, d),
3   (profit1, profit2) -> (profit1 + profit2)
4 end

```

or maybe everybody pirates his music and David gets paid nothing:

```

Snippet 10
1 function €clone(dprice, d)
2   (0, 0), (d, d),
3   (profit1, profit2) -> 0
4 end

```

or maybe the participants of our economy are playing hot potato, and always exactly one of the participants needs to pay:

```

Snippet 11
1 function €clone(dprice, d)
2   rand(Bool) ? (dprice, 0) : (0, dprice), (d, d),
3   (profit1, profit2) -> (profit1 + profit2)
4 end

```

However, note that while this is a valid definition under our definition of an optic functor, applying € now leads to accumulation order dependence (the same happens in the variant where cloning is done once per value). As a result, € would no longer preserve standard SSA invariants. This is legal according to our definition, but it can be convenient to be able to arbitrarily permute SSA transforms and optic functors. Thus, we would generally only ever choose one of the first two definitions.

²Welcome to Copyright 101 - Today's topic: Category theory.

2) *Delete*: Like copy, we have some freedom what do here. For our example, it would probably make sense to dynamically disallow deleting objects that have a non-zero price, so we might write:

```

Snippet 12
1 function €delete(dprice, d)
2   dprice != 0 && error("Who's gonna pay for that?")
3   () -> (0, )
4 end

```

3) *Branch*: Again we have several potential choices. The primary reason that the branch operation was non-canonical was because of the implicit delete of the branch condition. However, it is possible for us to implement more complicated behaviors. For example suppose we wanted to proportionately assign to each input. We might write something like (encoding for the branch functor is slightly tricky, this is one example):

```

Snippet 13
1 function €branch(aprices, cprice, as, c)
2   # Smear the cost of 'c' proportionately
3   total_a = sum(aprices)
4   aprices = map(p -> p * (1 + cprice / total_a), aprices)
5   aprices, (c, as), profit -> begin
6     (profit - cprice), cprice
7   end
8 end

```

However, as with the prior discussion, this would not match standard ssa behavior (i.e. programs equivalent under SSA would not longer be equivalent when lifted). As before, this is legal, but may not be desirable.

IV. TANGENT SPACES

We will briefly define the concepts relevant to us, but otherwise defer to any introductory text on differential geometry³. Let M be a smooth manifold, $C^\infty(M)$ the space of smooth functions on M . A tangent vector ν to M at a is an \mathbb{R} -linear derivation from $C^\infty(M)$ into \mathbb{R} , i.e. a function from $C^\infty(M) \rightarrow \mathbb{R}$ satisfying:

$$\nu(f \cdot g)(x) = g(a) \cdot \nu(f) + f(a) \cdot \nu(g) \quad (12)$$

The collection of all tangent vectors to M at a form a vector space, called the tangent space to M at a , $T_a M$. Additionally, if we glue all these vector spaces together, we obtain the tangent bundle TM .

We can dualize this vector space and obtain the cotangent bundle T^*M . Sections of the cotangent bundle are also called differential 1-forms.

If we have a function between manifolds, say $f : M \rightarrow N$, we obtain a pushforward of tangent spaces: $f_* : T_M \rightarrow T_N$ given by:

$$(f_* \nu)g = \nu(g \circ f) \quad (13)$$

(for $g : C^\infty(M)$), as well as a pullback of cotangent spaces $f^* : T^*N \rightarrow T^*M$ given by

$$(f^* dx)\nu = dx(f_* \nu) \quad (14)$$

³e.g. <https://faculty.math.illinois.edu/~lerman/518/f11/8-19-11.pdf>

For completeness, let us repeat this diagram, making use of the notation we established above for optic obtained from optic functors:



For completeness, let us proceed for one more step and explicitly draw a 3-optic. As before, we have a pair of two 2-optics:



Now, writing down the 3-optic is simply a matter of matching up the M s, and setting corresponding gaps opposite each other. This is done in figure 2. Additionally, the composition rule for 3 optics is shown in the appendix in figure 4.

C. Well-behaved primitives

We say a primitive of an optic functor ϵ is *well behaved* if the LHS of ϵf is simply the original function 'f', possibly with a capture of the input or output value, i.e. $\epsilon f = \langle \Delta_{SSA} \circ (id \otimes f) | f^* \rangle$ or $\epsilon f = \langle f \circ \Delta_{SSA} | f^* \rangle$. For functions of this form, various higher order concepts simplify and we can perform analysis efficiently using dynamic programming. Specifically, for an order n transform, we only need to do $n^2 - 1$ work rather than 2^n as would naively be expected. To see this, note that by repeatedly applying the transform, we can commute any function in our n -optic to just be a structural transform together with $f^{**}\dots$ where the number of stars is the total number of stars in the original expression. Then, letting numbers indicate the number of stars, we get the following sequence of patterns for increasing optics: 0, 01, 0121, 01212321, 0121232123432321, which is amenable to dynamic programming analysis.

Is this actually a useful notion? Certainly more advanced. Let's leave it here for now.

VI. REVERSE MODE AD

There are many ways to describe the difference between forward and reverse mode AD, but using our newfound language, it becomes quite simple: Forward mode AD pushes forward tangent vectors, while reverse mode AD pulls back cotangent vectors.

Now for the payoff from all of our hard work from the previous sections: Computing pullbacks of differential forms has the structure of an optic. As a concrete example, suppose we wanted to compute the gradient ∇f . Tracing through our definitions above, we find that:

$$\nabla f(x) = f^*(dx_1(f(x))) \quad (24)$$

so to evaluate this function we first have a covariant evaluation of f , followed by a contravariant evaluation of f^* . Now for the formal construction: Let \mathcal{C} be the category of Riemannian manifolds. We have an optic category $\mathbf{Optic}_{\mathcal{C}}(T^*M)$ ⁴ that pairs to each manifold its tangent bundle. Now, guided by the above, we defined a functor:

$$\overleftarrow{\partial} f : \mathcal{C} \rightarrow \mathbf{Optic}_{\mathcal{C}}(T^*M) \quad (25)$$

where as indicated, $\overleftarrow{\partial} M = T^*M$ and on morphisms $\overleftarrow{\partial} f = \langle (a \rightarrow (a, f(a)), (a, z) \rightarrow f^*_a(z)) \rangle$.

We call this functor the **Reverse mode AD optic functor**.

There is very little left to do at this point, but our program above did teach us that we need to define a few things:

For copy, we simply write

$$\overleftarrow{\partial} copy(x) = (x, x), (\Delta_1, \Delta_2) \rightarrow \Delta_1 + \Delta_2 \quad (26)$$

which is the unique map that preserves linear of the pullback. For delete:

$$\overleftarrow{\partial} delete(x) = () \rightarrow 0 \quad (27)$$

For branch we have a choice. We could either inherit the behavior from delete or return a poison value on the backwards pass. This is a user level policy decision.

Lastly, we of course need to define derivatives for our primitives. We defer this task to ChainRules.jl.

VII. HIGHER ORDER AD

A. Straightforwardly

Let us investigate the structure of higher order AD. One straightforward way to obtain the higher order transform is to simply treat the encoded optic as yet another SSA program. As an example, we may explicitly do the transform $\overleftarrow{\partial} \overleftarrow{\partial} \sin$:

Snippet 14

```

1 function  $\overleftarrow{\partial} \overleftarrow{\partial} \sin(x)$ 
2    $\sin(x), \Delta \rightarrow$ begin
3      $\Delta * \cos(x), \Delta' \rightarrow$ begin
4        $\Delta' * \Delta * (-\sin(x)), \Delta' * \cos(x)$ 
5     end
6   end,  $(\Delta''', \Delta'') \rightarrow$ begin
7      $\Delta''' * \cos(x) + \Delta''$ 
```

⁴We're cheating here a little bit by calling the pair (A, A') simply T^*M , technically we have a dependent optic $(x : M, T_x M)$, but since our programming language does not care about such things, we're happy to cheat and let others work out the tedious details.



Fig. 1. Composition of two second order optics layed out horizontally. Indicated by dashed lines are the cut/glue lines of the original optic (as well as the resulting optic, which corresponds with the cutline of the second optic).



Fig. 2. Wiring diagram for the canonical 3-optic $(A_{1-8}) \rightarrow (B_{1-8})$

8 **end**
9 **end**

Recall from above the diagram for an optic functor

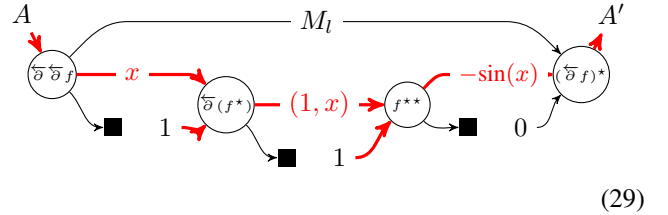


(28)

We can identify the various arrows with concrete values in the program:

Diagram	Program
M_l	Capture of x in the last closure
$M_{\overleftarrow{\partial}}$	Capture of x in the first closure
B	Return value of $\overleftarrow{\partial} \overleftarrow{\partial} \sin$
B'	Δ
M_r	Capture of Δ, x by the inner closure
A'	Primal return value of the first closure ($\Delta * \cos(x)$)
A''	Δ'
B''	Second return, inner closure ($\Delta' * \cos(x)$)
M'_l	First return, inner closure ($-\Delta' * \Delta * \sin(x)$)
B'	Also, Δ''' Δ''

In particular, if we are taking the second derivative of sin, here is the dataflow path of the final answer:



(29)

For a more complex example, the same transform is done for a third order transform in the Appendix (figure 6).

B. Alternative encoding

The previous approach works, but it is somewhat unsatisfying. Why? For me the two primary reasons are that there is no manifest difference between the residuals and the “holes” of the optic and that control flow is somewhat non-transparent. At first order, the residual was encoded as the capture of the original closure. However, at second order, while most residuals are still closure captures, Δ'' is a residual, but does not look like a capture. This shows up more at third order. To see what I mean when I say control-flow is non-transparent, consider the order in which the functions are called at third order:

$$\overleftarrow{\partial} \overleftarrow{\partial} f, \overleftarrow{\partial} \overleftarrow{\partial} (f^*), \overleftarrow{\partial} (f^{**}), \overleftarrow{\partial} ((\overleftarrow{\partial} f)^*), (\overleftarrow{\partial} f)^{**}, f^{***}, (\overleftarrow{\partial} (f^*))^*, (\overleftarrow{\partial} \overleftarrow{\partial} f)^* \quad (30)$$

(if it's not clear how to come up with this order, simply look at the second order diagram, add $\overleftarrow{\partial}$ in front of everything, then reverse the diagram, add $*$ after everything

and append the reversed diagram to the one that had $\overleftarrow{\partial}$ added). Comparing this to figure 6 in the appendix, we see that in source code order, the sequence is 1, 2, 3, 6, 7, 5, 8, so we need to do a lot of jumping around to keep track of what's happening. Can we re-arrange the encoding to make it easier to understand what's going on? Let's take a clue from the diagrams we've been working with, which do make the control flow order manifestly clear. What if instead of our second order example above, we wrote it as:

```

1  function  $\partial^2 \sin(x)$ 
2      sin(x),  $\Delta \rightarrow$ begin
3           $\Delta * \cos(x)$ ,  $\Delta' \rightarrow$ begin
4               $\Delta' * \cos(x)$ , let  $\Delta''' = \Delta' * \Delta * (-\sin(x))$ 
5                   $\Delta'' \rightarrow \Delta'' * \cos(x) + \Delta'''$ 
6              end
9          end
10     end
11 end

```

This avoids the two problems identified above: The Δ''' becomes an explicit capture and the order of operations becomes quite clear (outer functions are evaluated before their inner closures). This representation is a bit of a re-think of the standard transform, but it is entirely equivalent. To get more familiar with it, let's say, we'd like to take the third order transform of the following functions:

```

1  function f(x)
2      y = g(x)
3      z = h(y)
4      return z
5  end

```

We would have:

```

1  function  $\overleftarrow{\partial}^3 f(x)$ 
2      y,  $\overline{y} = \overleftarrow{\partial}^3 g(x)$ 
3      z,  $\overline{z} = \overleftarrow{\partial}^3 h(y)$ 
4      z,  $\Delta \rightarrow$ begin
5           $\alpha$ ,  $\overline{\alpha} = \overline{z}(\Delta)$ 
6           $\beta$ ,  $\overline{\beta} = \overline{y}(\alpha)$ 
7           $\beta$ ,  $\Delta' \rightarrow$ begin
8               $\gamma$ ,  $\overline{\gamma} = \overline{\beta}(\Delta')$ 
9               $\delta$ ,  $\overline{\delta} = \overline{\alpha}(\gamma)$ 
10              $\delta$ ,  $\Delta'' \rightarrow$ begin
11                  $\epsilon$ ,  $\overline{\epsilon} = \overline{\delta}(\Delta'')$ 
12                  $\zeta$ ,  $\overline{\zeta} = \overline{\gamma}(\epsilon)$ 
13                  $\zeta$ ,  $\Delta''' \rightarrow$ begin
14                      $\eta$ ,  $\overline{\eta} = \overline{\zeta}(\Delta''')$ 
15                      $\theta$ ,  $\overline{\theta} = \overline{\epsilon}(\eta)$ 
16                      $\theta$ ,  $\Delta^4 \rightarrow$ begin
17                          $\iota$ ,  $\overline{\iota} = \overline{\theta}(\Delta^4)$ 
18                          $\kappa$ ,  $\overline{\kappa} = \overline{\eta}(\iota)$ 
19                          $\kappa$ ,  $\Delta^5 \rightarrow$ begin
20                              $\lambda$ ,  $\overline{\lambda} = \overline{\kappa}(\Delta^5)$ 
21                              $\mu$ ,  $\overline{\mu} = \overline{\iota}(\lambda)$ 
22                              $\mu$ ,  $\Delta^6 \rightarrow$ begin
23                                  $\nu$ ,  $\overline{\nu} = \overline{\mu}(\Delta^6)$ 
24                                  $\xi$ ,  $\overline{\xi} = \overline{\lambda}(\nu)$ 
25                                  $\xi$ ,  $\Delta^7 \rightarrow$ begin
26                                      $o = \overline{\xi}(\Delta^7)$ 
27                                      $\pi = \overline{\nu}(o)$ 
28                                 return  $\pi$ 
29                             end
30                         end
31                     end
32                 end
33             end
34         end
35     end
36 end

```

```

1  end
2  end
3  end
4  end

```

We see still have the same number of functions, so the complexity isn't reduced, but at least it's fairly orderly. We do have to take special care of primitives, because we're no longer just applying the same transformation to the rule function itself. So suppose we have some primitive rule like:

```

1  struct  $\nabla a$ 
2      x
3  end
4  (:: $\nabla a$ ) = ...
5
6  function rrule(::typeof(a), x)
7      b(x),  $\nabla a(x)$ 
8  end

```

For the first order, we would just expose the raw pullback as we do with the naive transform. However, for second order, things become slightly more complicated:

```

1  function  $\overleftarrow{\partial}^2 a(x)$ 
2      (y,  $\overline{y}$ ),  $\overline{\overline{y}} = \overleftarrow{\partial}^2 \text{rrule}(a, x)$ 
3      y,  $\Delta \rightarrow$ begin
4          ( $\alpha_1$ ,  $\alpha_2$ ),  $\overline{\alpha} = \overleftarrow{\partial}^2 \text{typeof}(\overline{y})(\overline{y}, \Delta)$ 
5          ( $\alpha_1$ ,  $\alpha_2$ ), ( $\Delta'_1$ ,  $\Delta'_2$ )  $\rightarrow$ begin
6              ( $\Delta'''$ ,  $\beta$ ) =  $\overline{\alpha}((\Delta'_1, \Delta'_2))$ 
7               $\beta$ ,  $\Delta'' \rightarrow$ begin
8                  # Drop gradient w.r.t. `rrule`
9                  ( $\_$ ,  $a'$ ,  $x'$ ) =  $\overline{\overline{y}}(\Delta'', \Delta''')$ 
10                 return ( $a'$ ,  $x'$ )
11             end
12         end
13     end
14 end

```

It's worth comparing this with the explicitly written out and inlined version of sin above (which is a primitive).

It's also worth taking a look at the third order version of that example to see the generalizing pattern:

```

1  function  $\overleftarrow{\partial}^3 a(x)$ 
2      (y,  $\overline{y}$ ),  $\overline{\overline{y}} = \overleftarrow{\partial}^3 \text{rrule}(a, x)$ 
3      y,  $\Delta \rightarrow$ begin
4          (( $\alpha_1$ ,  $\alpha_2$ ),  $\overline{\alpha}$ ) =  $\overleftarrow{\partial}^3 \text{typeof}(\overline{y})(\overline{y}, \Delta)$ 
5          ( $\alpha_1$ ,  $\alpha_2$ ), ( $\Delta'_1$ ,  $\Delta'_2$ )  $\rightarrow$ begin
6              ( $\Delta'''$ ,  $\beta$ ),  $\overline{\beta} = \overline{\alpha}((\Delta'_1, \Delta'_2))$ 
7               $\beta$ ,  $\Delta'' \rightarrow$ begin
8                  ( $\_$ ,  $\gamma_1$ ,  $\gamma_2$ ),  $\overline{\gamma} = \overline{\overline{y}}(\Delta'', \Delta''')$ 
9                  ( $\gamma_1$ ,  $\gamma_2$ ), ( $\Delta^4_1$ ,  $\Delta^4_2$ )  $\rightarrow$ begin
10                     ( $\delta_1$ ,  $\delta_2$ ),  $\overline{\delta} = \overline{\gamma}((\Delta^4_1, \Delta^4_2))$ 
11                      $\delta_1$ ,  $\Delta^5 \rightarrow$ begin
12                         ( $\epsilon_1$ ,  $\epsilon_2$ ),  $\overline{\epsilon} = \overline{\overline{\beta}}(\delta_2, \Delta^5)$ 
13                         ( $\epsilon_1$ ,  $\epsilon_2$ ), ( $\Delta^6_1$ ,  $\Delta^6_2$ )  $\rightarrow$ begin
14                             ( $\zeta_1$ ,  $\zeta_2$ ) =  $\overline{\overline{\epsilon}}((\Delta^6_1, \Delta^6_2))$ 
15                              $\zeta_2$ ,  $\Delta^7 \rightarrow$ begin
16                                 return  $\overline{\delta}(\Delta^7, \zeta_1)$ 
17                             end
18                         end
19                     end
20                 end
21             end
22         end
23     end
24 end

```

If this structure seems familiar, it should! In our 3-optic

diagram we build 3-optics as pairs of 2-optics that composed in a canonical way. This is just the very same thing (compare functions l_{x_x} with closures returned from $\overleftarrow{\partial}^2 \text{rrule}$ and r_{x_x} with closures returned from $\overleftarrow{\partial}^2 \text{typeof}(\overline{y})$).

There is one additional case that's worth considering: We allow explicitly calling the derivative operator in the primal. What happens when we encounter such a function. The answer is fairly straightforward. We're basically doing the inverse of the rrule case above. In particular, we're adapting an $n + 1$ optic to an n optic by appropriate interleaving. Suppose we have a primal function ∇ given by:

Snippet 21

```

1 function  $\nabla(f, x)$ 
2    $(a, \overline{a}) = \overleftarrow{\partial}(f, x)$ 
3    $y = g(a)$ 
4    $z = \overline{a}(y)$ 
5   return  $z$ 
6 end

```

Applying our transform, we get:

Snippet 22

```

1 function  $\overleftarrow{\partial}\nabla(f, x)$ 
2    $(a, \overline{a}) = \overleftarrow{\partial}^2(f, x)$ 
3    $y, \overline{y} = \overleftarrow{\partial}g(a)$ 
4    $z, \overline{z} = \overline{a}(y)$ 
5    $z, \Delta \rightarrow \text{begin}$ 
6      $\alpha, \overline{\alpha} = \overline{z}(\Delta)$ 
7      $\beta = \overline{y}(\alpha)$ 
8      $\epsilon = \overline{\alpha}(\beta)$ 
9     return  $\epsilon$ 
10  end
11 end

```

Note that we went through all four callbacks of the 2-optic $\overleftarrow{\partial}^2(f, x)$, even though the top level optic only had two callbacks. For completeness, here is the third order version (second order transform):

Snippet 23

```

1 function  $\overleftarrow{\partial}^2\nabla(f, x)$ 
2    $(a, \overline{a}) = \overleftarrow{\partial}^3(f, x)$ 
3    $y, \overline{y} = \overleftarrow{\partial}g(a)$ 
4    $z, \overline{z} = \overline{a}(y)$ 
5    $z, \Delta \rightarrow \text{begin}$ 
6      $\alpha, \overline{\alpha} = \overline{z}(\Delta)$ 
7      $\beta, \overline{\beta} = \overline{y}(\alpha)$ 
8      $\epsilon, \overline{\epsilon} = \overline{\alpha}(\beta)$ 
9      $\epsilon, \Delta' \rightarrow \text{begin}$ 
10        $\lambda, \overline{\lambda} = \overline{\epsilon}(\Delta')$ 
11        $\zeta, \overline{\zeta} = \overline{\beta}(\lambda)$ 
12        $\mu, \overline{\mu} = \overline{\lambda}(\zeta)$ 
13        $\mu, \Delta'' \rightarrow \text{begin}$ 
14          $\theta, \overline{\theta} = \overline{\mu}(\Delta'')$ 
15          $\iota = \overline{\zeta}(\theta)$ 
16          $\kappa = \overline{\theta}(\iota)$ 
17         return  $\kappa$ 
18       end
19     end
20  end
21 end

```

VIII. IMPLEMENTATION DETAILS

The key to a performant reverse mode AD implementation is to pick efficient residual representatives M for our optic as well as introducing minimal additional overhead in the tensor. In julia, the latter means being amenable to type inference (as well as ideally not being significantly

harder on type inference than a simple primal inference to avoid unexpected performance cliffs). So, how do perform type inference on morphisms from the optic? Well, if we pick a particular representative, we obtain morphisms in the underlying category that we can do with whatever we want, including performing inference. However, picking a representative of course fixes the residual, which we'd like to avoid. Additionally, we run into problems with nested application of the functor.

IX. APPENDIX

Printing the appendix is encouraged to use as a template for working through examples.

REFERENCES

- [1] M. Riley, “Categories of optics,” 2018. [Online]. Available: <https://arxiv.org/abs/1809.00738>



Fig. 3. Composition rule for 2-optics (larger version)



Fig. 4. Composition rule for 3-optics

```

1  function  $\overleftarrow{\partial} \overleftarrow{\partial} \sin(x)$ 
2       $\sin(x), \Delta \rightarrow$  begin
3           $\Delta * \cos(x), \Delta' \rightarrow$  begin
4               $\Delta' * \Delta * (-\sin(x)), \Delta' * \cos(x)$ 
5          end
6      end,  $(\Delta''', \Delta'') \rightarrow$  begin
7           $\Delta'' * \cos(x) + \Delta'''$ 
8      end
9  end

```

Fig. 5. Explicitly written 2-nd order example (larger version)

```

1  function  $\overleftarrow{\partial} \overleftarrow{\partial} \overleftarrow{\partial} \sin(x)$ 
2       $(\sin(x), \Delta \rightarrow$  begin #  $\overleftarrow{\partial} \overleftarrow{\partial} (f^*)$ 
3           $\Delta * \cos(x), \Delta' \rightarrow$  begin #  $\overleftarrow{\partial} (f^{**})$ 
4               $\Delta' * \Delta * (-\sin(x)), \Delta' * \cos(x), ((\Delta^4, \Delta^5),) \rightarrow$  begin #  $f^{***}$ 
5                   $(-\Delta^4 * \Delta' * \Delta * \cos(x) +$ 
6                       $-\Delta^5 * \Delta' * \sin(x),$ 
7                       $-\Delta^4 * \Delta' * \sin(x)),$ 
8                       $-\Delta^4 * \Delta * \sin(x) + \Delta^5 * \cos(x)$ 
9              end
10             end,  $(\Delta^6, (\Delta^7, \Delta^8)) \rightarrow$  begin #  $(\overleftarrow{\partial} (f^*))^*$ 
11                  $(\Delta^6 * \Delta * (-\sin(x)) + \Delta^7,$ 
12                  $\Delta^6 * \cos(x) + \Delta^8)$ 
13             end
14         end,  $(\Delta''', \Delta'') \rightarrow$  begin #  $\overleftarrow{\partial} ((\overleftarrow{\partial} f)^*)$ 
15              $\Delta'' * \cos(x) + \Delta''', \Delta^9 \rightarrow$  begin #  $(\overleftarrow{\partial} f)^{**}$ 
16                  $(-\Delta^9 * \Delta'' * \sin(x)),$ 
17                  $(\Delta^9, \Delta^9 * \cos(x))$ 
18             end
19         end), (((\Delta^{10}, \Delta^{11}), \Delta^{12}),) \rightarrow begin #  $(\overleftarrow{\partial} \overleftarrow{\partial} f)^*$ 
20              $\Delta^{10} * \cos(x) + \Delta^{11} + \Delta^{12}$ 
21         end
22     end

```

Fig. 6. Explicitly written out third-order example