

Automatic differentiation in Coconut

Tom Ellis
Simon Peyton Jones
Andrew Fitzgibbon

Atoms		
f, g, h	$::=$	Function
x, y, z	$::=$	Local variable (lambda-bound or let-bound)
k	$::=$	Literal constants
Terms		
pgm	$::=$	$def_1 \dots def_n$
def	$::=$	$f(x) = e$
e	$::=$	k Constant
	$ $	x Local variable
	$ $	$f(e)$ Function call
	$ $	(e_1, e_2) Pair
	$ $	$\lambda x. e$ Lambda
	$ $	$e_1 e_2$ Application
	$ $	$let x = e_1 in e_2$
	$ $	$if b then e_1 else e_2$
Types		
τ	$::=$	\mathbb{N} Natural numbers
	$ $	\mathbb{R} Real numbers
	$ $	(τ_1, τ_2) Pairs
	$ $	$Vec \tau$ Vectors
	$ $	$\tau_1 \rightarrow \tau_2$ Functions
	$ $	$\tau_1 \multimap \tau_2$ Linear maps

Figure 1. Syntax of the language

1 The language

This paper is about automatic differentiation of functions, so we must be precise about the language in which those functions are written.

The syntax of our language is given in Figure 1. Note that

- Variables are divided into *functions*, f, g, h ; and *local variables*, x, y, z , which are either function arguments or let-bound.
- The language has a first order sub-language. Functions are defined at top level; functions always appear in a call, never (say) as an argument to a function; in a call $f(e)$, the function f is always a top-level-defined function, never a local variable. **AWF: at some point we should say where this restriction is needed**
- Functions have exactly one argument. If you want more than one, pass a pair.
- Pairs are built-in, with selectors $\pi_{1,2}, \pi_{2,2}$. In the real implementation, pairs are generalised to n -tuples, and we often do so informally here.
- Conditionals are a language construct. **SPJ: Treating “if” as a function just didn’t work; in particular ∇if needed a linear-map version of “if” and once we have that we might as well build “if” in. Anyway, conditionals are very fundamental, so it’s unsurprising.**
- Let-bindings are non-recursive. For now, at least, top-level functions are also non-recursive. **SPJ: I think that top-level recursive functions might be OK, but I don’t want to think about that yet.**
- Lambda expressions and applications are present, so the language is higher order. AD will only accept a subset of the language, in which lambdas appear only as an argument to *build*. But the *output* of AD may include lambdas and application, as we shall see.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

Conference’17, July 2017, Washington, DC, USA

© 2019 Association for Computing Machinery.

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM...\$15.00

<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

Built-in functions		
$(+)$	$:: \text{Field } t \Rightarrow (t, t) \rightarrow t$	
$(*)$	$:: \text{Field } t \Rightarrow (t, t) \rightarrow t$	
$\pi_{1,2}$	$:: (t_1, t_2) \rightarrow t_1$	Selection
$\pi_{2,2}$	$:: (t_1, t_2) \rightarrow t_2$..ditto..
build	$:: (\mathbb{N}, \mathbb{N} \rightarrow t) \rightarrow \text{Vec } t$	Vector build
ixR	$:: (\mathbb{N}, \text{Vec } t) \rightarrow t$	Indexing (NB arg order)
sum	$:: \text{Field } t \Rightarrow \text{Vec } t \rightarrow t$	Sum a vector
size	$:: \text{Vec } t \rightarrow \mathbb{N}$	Size of a vector
Derivatives of built-in functions		
$\nabla+$	$:: \text{Field } t \Rightarrow (t, t) \rightarrow ((t, t) \multimap t)$	
$\nabla+(x, y)$	$= 1 \bowtie 1$	
$\nabla*$	$:: \text{Field } t \Rightarrow (t, t) \rightarrow ((t, t) \multimap t)$	
$\nabla*(x, y)$	$= S(y) \bowtie S(x)$	
$\nabla\pi_{1,2}$	$:: (t, t) \rightarrow ((t, t) \multimap t)$	
$\nabla\pi_{1,2}(x)$	$= 1 \bowtie 0$	
∇ixR	$:: (\mathbb{N}, \text{Vec } t) \rightarrow ((\mathbb{N}, \text{Vec } t) \multimap t)$	
$\nabla\text{ixR}(i, v)$	$= 0 \bowtie \mathcal{B}'(\text{size}(v), \lambda j. \text{if } i = j \text{ then } 1 \text{ else } 0)$	
∇sum	$:: \text{Field } t \Rightarrow \text{Vec } t \rightarrow (\text{Vec } t \multimap t)$	
$\nabla\text{sum}(v)$	$= \mathcal{B}'(\text{size}(v), \lambda i. 1)$	
\dots		

Figure 2. Built-in functions

1.1 Built in functions

The language has built-in functions shown in Figure 2.

We allow ourselves to write functions infix where it is convenient. Thus $e_1 + e_2$ means the call $+(e_1, e_2)$, which applies the function $+$ to the pair (e_1, e_2) . (So, like all other functions, $+$ has one argument.) Similarly the linear map $m_1 \times m_2$ is short for $\times(e_1, e_2)$.

We allow ourselves to write vector indexing $\text{ixR}(i, a)$ using square brackets, thus $a[i]$.

Multiplication and addition are overloaded to work on any suitable type. On vectors they work element-wise; if you want dot-product you have to program it.

1.2 Vectors

The language supports one-dimensional vectors, of type $\text{Vec } T$, whose elements have type T (Figure 1). A matrix can be represented as a vector of vectors.

Vectors are supported by the following built-in functions (Figure 2):

- $\text{build} :: (\mathbb{N}, \mathbb{N} \rightarrow t) \rightarrow \text{Vec } t$ for vector construction.

- $\text{ixR} :: (\mathbb{N}, \text{Vec } t) \rightarrow t$ for indexing. Informally we allow ourselves to write $v[i]$ instead of $\text{ixR}(i, v)$.

- $\text{sum} :: \text{Field } t \Rightarrow \text{Vec } t \rightarrow t$ to add up the elements of a vector. We specifically do not have a general, higher order, fold operator; we say why in Section 3.4. **TE:** I believe that for a vector v of size n , $\text{sum}(v)$ is the same as $\mathcal{B}'(n, \text{const id}) v$. This may or may not be useful in reducing the size of the base language, should we want to do that. **SPJ:** I don't think so! \mathcal{B}' is a linear map, so you can't apply it to v . Maybe you mean $\mathcal{B}'(n, \text{const id}) \odot v$? But that (Figure 4) is defined using sum !

- $\text{size} :: \text{Vec } t \rightarrow \mathbb{N}$ takes the size of a vector.

- Arithmetic functions $(*)$, $(+)$ etc are overloaded to work over vectors, always elementwise.

SPJ: Do we need scan? Or (specialising to $(+)$) cumulative sum?

2 Linear maps and differentiation

If $f : S \rightarrow T$, then its derivative ∇f has type

$$\nabla f : S \rightarrow (S \multimap T)$$

where $S \multimap T$ is the type of *linear maps* from S to T . That is, at some point $p : S$, $\nabla f(p)$ is a linear map that is a good approximation of f at p .

By “a good approximation of f at p ” we mean this:

$$\forall p : S. f(p + \delta_p) \approx f(p) + \nabla f(p) \odot \delta_p$$

Here the operation \odot is linear-map application: it takes a linear map $S \multimap T$ and applies it to an argument of type S , giving a result of type T (Figure 3).

The linear maps from S to T are a subset of the functions from S to T . We characterise linear maps more precisely in Section 2.1, but a good intuition can be had for functions $g : \mathbb{R}^2 \rightarrow \mathbb{R}$. This function defines a curvy surface $z = g(x, y)$. Then a linear map of type $\mathbb{R}^2 \multimap \mathbb{R}$ is a plane, and $\nabla g(p_x, p_y)$ is the plane that best approximates g near (p_x, p_y) , that is a tangent plane passing through $z = g(p_x, p_y)$.

2.1 Linear maps

A *linear map*, $m : S \multimap T$, is a function from S to T , satisfying these two properties:

$$(LM1) \quad \forall x, y : S \quad m \odot (x + y) = m \odot x + m \odot y$$

$$(LM2) \quad \forall k : \mathbb{R}, x : S \quad k * (m \odot x) = m \odot (k * x)$$

Here $\odot : (s \multimap t) \rightarrow (s \rightarrow t)$ is an operator that applies a linear map $(s \multimap t)$ to an argument of type s .

The type $s \multimap t$ is a type in the language (Figure 1).

Linear maps can be *built and consumed* using the operators in (see Figure 3). Indeed, you should think of linear maps as an *abstract type*; that is, you can *only* build or consume linear maps with the operators in Figure 3. We might *represent* a

Operator	Type	Matrix interpretation
where $s = \mathbb{R}^m$, and $t = \mathbb{R}^n$		
Apply $(\odot) : (s \multimap t) \rightarrow (s \rightarrow t)$		Matrix/vector multiplication
Compose $(\circ) : (s \multimap t, r \multimap s) \rightarrow (r \multimap t)$		Matrix/matrix multiplication
Sum $(\oplus) : \text{Field } t \Rightarrow (s \multimap t, s \multimap t) \rightarrow (s \multimap t)$		Matrix addition
Zero $\mathbf{0} : \text{Field } t \Rightarrow s \multimap t$		Zero matrix
Unit $\mathbf{1} : s \multimap s$		Identity matrix (square)
Scale $\mathcal{S}(\cdot) : \text{Field } s \Rightarrow s \rightarrow (s \multimap s)$		
Pair $(\times) : \text{Field } s \Rightarrow (s \multimap t_1, s \multimap t_2) \rightarrow (s \multimap (t_1, t_2))$		Vertical juxtaposition
Join $(\bowtie) : \text{Field } s \Rightarrow (t_1 \multimap s, t_2 \multimap s) \rightarrow ((t_1, t_2) \multimap s)$		Horizontal juxtaposition
Transpose $\cdot^\top : (s \multimap t) \rightarrow (t \multimap s)$		Matrix transpose
NB: We expect to have only \mathcal{L}/\mathcal{L}' or \mathcal{B}/\mathcal{B}', but not both		
Lambda $\mathcal{L} : (\mathbb{N} \rightarrow (s \multimap t)) \rightarrow (s \multimap (\mathbb{N} \rightarrow t))$		
TLambda $\mathcal{L}' : (\mathbb{N} \rightarrow (t \multimap s)) \rightarrow ((\mathbb{N} \rightarrow t) \multimap s)$		Transpose of \mathcal{L}
Build $\mathcal{B} : (\mathbb{N}, \mathbb{N} \rightarrow (s \multimap t)) \rightarrow (s \multimap \text{Vec } t)$		
BuildT $\mathcal{B}' : (\mathbb{N}, \mathbb{N} \rightarrow (t \multimap s)) \rightarrow (\text{Vec } t \multimap s)$		Transpose of \mathcal{B}

Figure 3. Operations over linear maps

linear map in a variety of ways, one of which is as a matrix (Section 2.4).

2.1.1 Semantics

The *semantics* of a linear map is completely specified by saying what ordinary function it corresponds to; or, equivalently, by how it behaves when applied to an argument by (\odot) . The semantics of each form of linear map are given in Figure 4

2.1.2 Laws

Linear maps satisfy *laws* given in Figure 4. Note that (\circ) and \oplus behave like multiplication and addition respectively.

These laws can readily be proved from the semantics. To prove two linear maps are equal, we must simply prove that they give the same result when applied to any argument. So, to prove that $\mathbf{0} \circ m = m$, we choose an arbitrary x and reason thus:

$$\begin{aligned}
 & (\mathbf{0} \circ m) \odot x \\
 &= \mathbf{0} \odot (m \odot x) \quad \{\text{semantics of } (\odot)\} \\
 &= \mathbf{0} \quad \{\text{semantics of } \mathbf{0}\} \\
 &= \mathbf{0} \odot x \quad \{\text{semantics of } \mathbf{0} \text{ backwards}\}
 \end{aligned}$$

Theorem: $\forall(m : S \multimap T). m \odot \mathbf{0} = \mathbf{0}$. That is, all linear maps pass through the origin. **Proof:** property (LM2) with $k = 0$. Note that the function $\lambda x.x + 4$ is not a linear map; its graph is a straight line, but it does not go through the origin.

2.2 Vector spaces

Given a linear map $m : S \multimap T$, we expect both S and T to be a *normed vector space*¹. A vector space V has:

- *Vector addition* $(+_V) : V \rightarrow V \rightarrow V$.
- *Zero vector* $0_V : V$.
- *Scalar multiplication* $(*_V) : \mathbb{R} \rightarrow V \rightarrow V$
- *Dot-product* $(\bullet_V) : V \rightarrow V \rightarrow \mathbb{R}$. This operation is what makes it a *normed* vector space.

We omit the V subscripts when it is clear which $(*)$, $(+)$, (\bullet) or 0 is intended.

These operations must obey the laws of vector spaces

$$\begin{aligned}
 v_1 + (v_2 + v_3) &= (v_1 + v_2) + v_3 \\
 v_1 + v_2 &= v_2 + v_1 \\
 v + 0 &= 0 \\
 0 * v &= 0 \\
 1 * v &= v \\
 r_1 * (r_2 * v) &= (r_1 * r_2) * v \\
 r * (v_1 + v_2) &= (r * v_1) + (r * v_2) \\
 (r_1 + r_2) * v &= (r_1 * v) + (r_2 * v)
 \end{aligned}$$

2.3 Transposition

For any linear map $m : S \multimap T$ we can produce its transpose $m^\top : T \multimap S$. Despite its suggestive type, the transpose is *not* the inverse of m ! (In the world of matrices, the transpose of a matrix is not the same as its inverse.)

¹https://en.wikipedia.org/wiki/Vector_space

Semantics of linear maps	
$(m_1 \circ m_2) \odot x$	$= m_1 \odot (m_2 \odot x)$
$(m_1 \times m_2) \odot x$	$= (m_1 \odot x, m_2 \odot x)$
$(m_1 \bowtie m_2) \odot (x_1, x_2)$	$= (m_1 \odot x_1) + (m_2 \odot x_2)$
$(m_1 \oplus m_2) \odot x$	$= (m_1 \odot x) + (m_2 \odot x)$
$0 \odot x$	$= 0$
$1 \odot x$	$= x$
$S(k) \odot x$	$= k * x$
$\mathcal{L}(f) \odot x$	$= \lambda i. (f \ i) \odot x$
$\mathcal{L}'(f) \odot g$	$= \Sigma_i (f \ i) \odot g(i)$
$\mathcal{B}(n, \lambda i. m) \odot x$	$= \text{build}(n, \lambda i. m \odot x)$
$\mathcal{B}'(n, \lambda i. m) \odot x$	$= \text{sum}(\text{build}(n, \lambda i. m \odot x[i]))$
Laws for linear maps	
$0 \circ m$	$= 0$
$m \circ 0$	$= 0$
$1 \circ m$	$= m$
$m \circ 1$	$= m$
$m \oplus 0$	$= m$
$0 \oplus m$	$= m$
$m \circ (n_1 \bowtie n_2)$	$= (m \circ n_1) \bowtie (m \circ n_2)$
$(m_1 \bowtie m_2) \circ (n_1 \times n_2)$	$= (m_1 \circ n_1) \oplus (m_2 \circ n_2)$
$S(k_1) \circ S(k_2)$	$= S(k_1 * k_2)$
$S(k_1) \oplus S(k_2)$	$= S(k_1 + k_2)$

Figure 4. Laws for linear maps

Laws for transposition of linear maps	
$(m_1 \circ m_2)^\top$	$= m_2^\top \circ m_1^\top$ Note reversed order!
$(m_1 \times m_2)^\top$	$= m_1^\top \bowtie m_2^\top$
$(m_1 \bowtie m_2)^\top$	$= m_1^\top \times m_2^\top$
$(m_1 \oplus m_2)^\top$	$= m_1^\top \oplus m_2^\top$
0^\top	$= 0$
1^\top	$= 1$
$S(k)^\top$	$= S(k)$
$(m^\top)^\top$	$= m$
$\mathcal{B}(n, \lambda i. m)^\top$	$= \mathcal{B}'(n, \lambda i. m^\top)$
$\mathcal{B}'(n, \lambda i. m)^\top$	$= \mathcal{B}(n, \lambda i. m^\top)$
$\mathcal{L}(\lambda i. m)^\top$	$= \mathcal{L}'(\lambda i. m^\top)$
$\mathcal{L}'(\lambda i. m)^\top$	$= \mathcal{L}(\lambda i. m^\top)$

Figure 5. Laws for transposition

Definition 2.1. Given a linear map $m : S \multimap T$, its *transpose* $m^\top : T \multimap S$ is defined by the following property:

$$(TP) \quad \forall s:S, t:T. (m^\top \odot t) \bullet s = t \bullet (m \odot s)$$

This property *uniquely* defines the transpose, as the following theorem shows:

Theorem 2.2. If m_1 and m_2 are linear maps satisfying

$$\forall s t. (m_1 \odot t) \bullet s = (m_2 \odot t) \bullet s$$

then $m_1 = m_2$

Proof. It is a property of dot-product that if $v_1 \bullet x = v_2 \bullet x$ for every x , then $v_1 = v_2$. (Just use a succession of one-hot vectors for x , to pick out successive components of v_1 and v_2 .) That means that (for every t):

$$m_1 \odot t = m_2 \odot t$$

and that is the definition of extensional equality. So m_1 and m_2 are the same linear maps. \square

Figure 5 has a collection of laws about transposition. These identities are readily proved using the above definition. For example, to prove that $(m_1 \circ m_2)^\top = m_2^\top \circ m_1^\top$ we may reason as follows:

$$\begin{aligned}
 & ((m_2^\top \circ m_1^\top) \odot t) \bullet s \\
 &= (m_2^\top \odot (m_1^\top \odot t)) \bullet s \quad \text{Semantics of } (\circ) \\
 &= (m_1^\top \odot t) \bullet (m_2 \odot s) \quad \text{Use (TP)} \\
 &= t \bullet (m_1 \odot (m_2 \odot s)) \quad \text{Use (TP) again} \\
 &= t \bullet ((m_1 \circ m_2) \odot s) \quad \text{Semantics of } (\circ)
 \end{aligned}$$

And now the property follows by Theorem 2.2.

2.4 Matrix interpretation of linear maps

A linear map $m : \mathbb{R}^M \multimap \mathbb{R}^N$ is isomorphic to a matrix $\mathbb{R}^{N \times M}$ with N rows and M columns.

Many of the operators over linear maps then have simple matrix interpretations; for example, composition of linear maps (\circ) is matrix multiplication, pairing (\times) is vertical juxtaposition, and so on. These matrix interpretations are all given in the final column of Figure 3.

You might like to check that matrix transposition satisfies property (TP).

When it comes to implementation, we do not want to *represent* a linear map by a matrix, because a linear map $\mathbb{R}^M \multimap \mathbb{R}^N$ is an $N \times M$ matrix, which is enormous if $N = M = 10^6$, say. The function might be very simple (even the identity function) and taking 10^{12} numbers to represent it is plain silly. So our goal will be to *avoid* realising linear maps as matrices.

2.5 Optimisation

In optimisation we are usually given a function $f : \mathbb{R}^N \rightarrow \mathbb{R}$, where N can be large, and asked to find values of the input that maximises the output. One way to do this is by *gradient*

descent: start with a point p , make a small change to $p + \delta_p$, and so on. From p we want to move in the direction of maximum slope. (How *far* to move in that direction is another matter — indeed no one knows — but we will concentrate on the *direction* in which to move.)

Suppose $\delta(i, N)$ is the one-hot N-vector with 1 in the i 'th position and zeros elsewhere. Then $\delta_p[i] = \nabla f(p) \odot \delta(i, N)$ describes how fast the output of f changes for a change in the i 'th input. The direction of maximum slope is just the vector

$$\delta_p = (\delta_p[1] \ \delta_p[2] \ \dots \ \delta_p[N])$$

How can we compute this vector? We can simply evaluate $\nabla f(p) \odot \delta(i, N)$ for each i . But that amounts to running f N times, which is bad if N is large (say 10^6).

Suppose that we somehow had access to $\nabla^T f$. Then we can use property (TP), setting $\delta_f = 1$ to get

$$\forall \delta_p. \nabla f(p) \odot \delta_p = (\nabla^T f(p) \odot 1) \bullet \delta_p$$

Then

$$\begin{aligned} \delta_p[i] &= \nabla f(p) \odot \delta(i, N) \\ &= (\nabla^T f(p) \odot 1) \bullet \delta(i, N) \\ &= (\nabla^T f(p) \odot 1)[i] \end{aligned}$$

That is $\delta_p[i]$ is the i 'th component of $\nabla^T f(p) \odot 1$, so $\delta_p = \nabla^T f(p) \odot 1$.

That is, $\nabla^T f(p) \odot 1$ is the N-vector of maximum slope, the direction in which to move if we want to do gradient descent starting at p . And *that* is why the transpose is important.

2.6 Lambdas and linear maps

Notice the similarity between the type of (\times) and the type of \mathcal{L} ; the latter is really just an infinite version of the latter. Their semantics in Figure 4 are equally closely related.

The transpositions of these two linear maps, (\bowtie) and \mathcal{L}' , are similarly related. *But*, there is a problem with the semantics of \mathcal{L}' :

$$\mathcal{L}'(f) \odot g = \Sigma_i (f \ i) \odot g(i)$$

This is an *infinite sum*, so there is something fishy about this as a semantics.

2.7 Questions about linear maps

- Do we need 1? After all $\mathcal{S}(1)$ does the same job. But asking if $k = 1$ is dodgy when k is a float. **AWF: No, perfectly fine to ask if a float is 1 — the nearby floats are far away, and there's no other float f such that $\mathcal{S}(f) = 1$. SPJ: For the purposes of this paper I prefer having 1 as well; unity plays such a key role!**
- Do these laws fully define linear maps?

Notes

- In practice we allow n-ary versions of $m \bowtie n$ and $m \times n$.

Original function	$f : S \rightarrow T$ $f(x) = e$
Full Jacobian	$\nabla f : S \rightarrow (S \multimap T)$ $\nabla f(x) = \text{let } \nabla x = 1 \text{ in } \nabla_S \llbracket e \rrbracket$
Transposed Jacobian	$\nabla^T f : S \rightarrow (T \multimap S)$ $\nabla^T f(x) = (\nabla f(x))^\top$
Forward derivative	$f' : (S, S) \rightarrow T$ $f'(x, dx) = \nabla f(x) \odot dx$
Reverse derivative	$f' : (S, T) \rightarrow S$ $f'(x, dr) = \nabla^T f(x) \odot dr$

Differentiation of an expression

If $e : T$ then $\nabla_S \llbracket e \rrbracket : S \multimap T$	
$\nabla_S \llbracket k \rrbracket$	$= 0$
$\nabla_S \llbracket x \rrbracket$	$= \nabla x$
$\nabla_S \llbracket f(e) \rrbracket$	$= \nabla f(e) \circ \nabla_S \llbracket e \rrbracket$
$\nabla_S \llbracket (e_1, e_2) \rrbracket$	$= \nabla_S \llbracket e_1 \rrbracket \times \nabla_S \llbracket e_2 \rrbracket$
$\nabla_S \llbracket \text{build}(e_n, \lambda i. e) \rrbracket$	$= \mathcal{B}(e_n, \lambda i. \nabla_S \llbracket e \rrbracket)$
$\nabla_S \llbracket \lambda i. e \rrbracket$	$= \mathcal{L}(\lambda i. \nabla_S \llbracket e \rrbracket)$
$\nabla_S \llbracket \text{let } x = e_1 \text{ in } e_2 \rrbracket$	$= \text{let } x = e_1 \text{ in}$ $\text{let } \nabla x = \nabla_S \llbracket e_1 \rrbracket \text{ in}$ $\nabla_S \llbracket e_2 \rrbracket$

Figure 6. Automatic differentiation

3 AD as a source-to-source transformation

To perform source-to-source AD of a function f , we follow the plan outlined in Figure 6. Specifically, starting with a function definition $f(x) = e$:

- Construct the full Jacobian ∇f , and transposed full Jacobian $\nabla^T f$, using the transformations in Figure 6².
- Optimise these two definitions, using the laws of linear maps in Figure 4.
- Construct the forward derivative f' and reverse derivative f' , as shown in Figure 6³.
- Optimise these two definitions, to eliminate all linear maps. Specifically:
 - Rather than *calling* ∇f (in, say, f'), instead *inline* it.
 - Similarly, for each local let-binding for a linear map, of form $\text{let } \nabla x = e \text{ in } b$, inline ∇x at each of its

² We consider ∇f and $\nabla^T f$ to be the names of two new functions. These names are derived from, but distinctd from f , rather like f' or f_1 in mathematics.

³ Again f' and f' are new names, derived from f

occurrences in b . This may duplicate e ; but ∇x is a function that may be applied (via \odot) to many different arguments, and we want to specialise it for each such call. (I think.)

- Optimise using the rules of (\odot) in Figure 4.
- Use standard Common Subexpression Elimination (CSE) to recover any lost sharing.

Note that

- The transformation is fully compositional; each function can be AD'd independently. For example, if a user-defined function f calls another user-defined function g , we construct ∇g as described; and then construct ∇f . The latter simply calls ∇g .
- The AD transformation is *partial*; that is, it does not work for every program. In particular, it fails when applied to a lambda, or an application; and, as we will see in Section 3.3, it requires that *build* appears applied to a lambda.
- We give the full Jacobian for some built-in functions in Figure 6, including for conditionals (∇if).

3.1 Forward and reverse AD

Consider

$$f(x) = p(q(r(x)))$$

Just running the algorithm above on f gives

$$\begin{aligned} f(x) &= p(q(r(x))) \\ \nabla f(x) &= \nabla p \circ (\nabla q \circ \nabla r) \\ f'(x, dx) &= (\nabla p \circ (\nabla q \circ \nabla r)) \odot dx \\ &= \nabla p \odot ((\nabla q \circ \nabla r) \odot dx) \\ &= \nabla p \odot (\nabla q \odot (\nabla r \odot dx)) \\ \nabla^T f(x) &= (\nabla^T r \circ \nabla^T q) \circ \nabla^T p \\ f'(x, dr) &= ((\nabla^T r \circ \nabla^T q) \circ \nabla^T p) \odot dr \\ &= (\nabla^T r \circ \nabla^T q) \odot (\nabla^T p \odot dr) \\ &= \nabla^T r \odot (\nabla^T q \odot (\nabla^T p \odot dr)) \end{aligned}$$

In “The essence of automatic differentiation” Conal says (Section 12)

The AD algorithm derived in Section 4 and generalized in Figure 6 can be thought of as a family of algorithms. For fully right-associated compositions, it becomes forward mode AD; for fully left-associated compositions, reverse-mode AD; and for all other associations, various mixed modes.

But the forward/reverse difference shows up quite differently here: it has nothing to do with *right-vs-left association*, and everything to do with *transposition*.

This is mysterious. Conal is not usually wrong. I would like to understand this better. TE: I was also puzzled by this.

Conal's claim is suspicious to me, but firstly it's very cool and secondly it's Conal, so I want it to be true and I still hope it is.

3.2 Avoiding duplication

We may want to ANF-ise before AD to avoid gratuitous duplication. E.g.

$$\begin{aligned} \nabla_S[\text{sqrt}(x + (y * z))] \\ &= \nabla_{\text{sqrt}}(x + (y * z)) \circ \nabla_S[x + (y * z)] \\ &= \nabla_{\text{sqrt}}(x + (y * z)) \circ \nabla_+(x, y * z) \\ &\quad \circ (\nabla_S[x] \times \nabla_S[y * z]) \\ &= \nabla_{\text{sqrt}}(x + (y * z)) \circ \nabla_+(x, y * z) \\ &\quad \circ (\nabla x \times (\nabla^*(y, z) \circ (\nabla y \times \nabla z))) \end{aligned}$$

Note the duplication of $y * z$ in the result. Of course, CSE may recover it.

TE: Yes, although when I say “AD” I mean something that is distinct from what I mean by “symbolic differentiation”. In particular by “AD” I mean something that preserves sharing in a way that symbolic differentiation doesn't. Perhaps between us we should pin down some terminology. SPJ: I don't understand this. Perhaps you can make it precise?

TE: Consider $\exp(\exp(x))$. I consider its “symbolic derivative” to be $\exp(\exp(x))\exp(x)$ and its “forward automatic derivative” to be $\text{let } y = \exp(x) \text{ in } \exp(y)y$. In other words, taking proper care of sharing is what makes AD AD and not just any old form of symbolic differentiation, in my personal nomenclature at least. Does that make it any clearer what I mean? AWF: For me, “AD” very specifically implies a second argument to the function. That's how you detect it's AD. I.e. $f'(x, dx) = \dots$ is forward mode, and $f'(x, df) = \dots$ is reverse mode. There's a lot of chat about what AD really is, and those who want to avoid such chat often now say “algorithmic differentiation”, to mean “all this stuff”. The real claim we want to explore is this: “Forward mode is good for functions with small inputs and large outputs, e.g. $\mathbb{R} \mapsto \mathbb{R}^n$, and reverse mode is for $\mathbb{R}^n \mapsto \mathbb{R}$ ”.

3.3 AD for vectors

Like other built-in functions, each built-in function for vectors has its full Jacobian versions, defined in Figure 2. You may enjoy checking that ∇sum and ∇ixR are correct!

For *build* there are two possible paths, and it's not yet clear which is best

Direct path. Figure 6 includes a rule for $\nabla_S[\text{build}(e_n, \lambda i.e)]$. But *build* is an exception! It is handled specially by the AD transformation in Figure 6; there is no ∇build . Moreover the AD transformation only works if the second argument of the *build* is a lambda, thus $\text{build}(e_n, \lambda i.e)$. I tried dealing with *build* and lambdas separately, but failed (see Section ??).

I did think about having a specialised linear map for indexing, rather than using \mathcal{B}' , but then I needed its transposition,

so just using \mathcal{B}' seemed more economical. On the other hand, with the functions as I have them, I need the grotesquely delicate optimisation rule

```
sum(build(n, λi. if i == ei then e else 0))
= let i = ei in b
if i ∉ ei
```

I hate this!

3.4 General folds

We have $sum :: Vec \mathbb{R} \rightarrow \mathbb{R}$. What is ∇sum ? One way to define its semantics is by applying it:

$$\nabla sum :: Vec \mathbb{R} \rightarrow (Vec \mathbb{R} \multimap \mathbb{R})$$

$$\nabla sum(v) \odot dv = sum(dv)$$

That is OK. But what about product, which multiplies all the elements of a vector together? If the vector had three elements we might have

$$\nabla product([x_1, x_2, x_3]) \odot [dx_1, dx_2, dx_3]$$

$$= (dx_1 * x_2 * x_3) + (dx_2 * x_1 * x_3) + (dx_3 * x_1 * x_2)$$

This looks very unattractive as the number of elements grows. Do we need to use product?

This gives the clue that taking the derivative of *fold* is not going to be easy, maybe infeasible! Much depends on the particular lambda it appears. So I have left out product, and made no attempt to do general folds.

4 Implementation

The implementation differs from this document as follows:

- Rather than pairs, the implementation supports n -ary tuples. Similarly the linear maps (\times) and \bowtie are n -ary.
- Functions definitions can take n arguments, thus

$$f(x, y, z) = e$$

This is treated as equivalent to

$$f(t) = \text{let } x = \pi_{1,3}(t)$$

$$y = \pi_{2,3}(t)$$

$$z = \pi_{3,3}(t)$$

$$\text{in } e$$

5 Demo

You can run the prototype by saying `ghci Main`.

The function `demo :: Def -> IO ()` runs the prototype on the function provided as example. Thus:

```
bash$ ghci Main
```

```
*Main> demo ex2
```

```
-----
Original definition
-----
```

```
fun f2(x)
= let { y = x * x }
  let { z = x + y }
  y * z
```

```
-----
Anf-ised original definition
-----
```

```
fun f2(x)
= let { y = x * x }
  let { z = x + y }
  y * z
```

```
-----
The full Jacobian (unoptimised)
-----
```

```
fun Df2(x)
= let { Dx = lmOne() }
  let { y = x * x }
  let { Dy = lmCompose(D*(x, x), lmVCat(Dx, Dx)) }
  let { z = x + y }
  let { Dz = lmCompose(D+(x, y), lmVCat(Dx, Dy)) }
  lmCompose(D*(y, z), lmVCat(Dy, Dz))
```

```
-----
The full Jacobian (optimised)
-----
```

```
fun Df2(x)
= let { y = x * x }
  lmScale( (x + y) * (x + x) + (x + y) * (x + x) )
```

```
-----
Forward derivative (unoptimised)
-----
```

```
fun f2'(x, dx)
= lmApply(let { y = x * x }
  lmScale( (x + y) * (x + x) +
           (x + y) * (x + x) ),
  dx)
```

```
-----
Forward-mode derivative (optimised)
-----
```

```
fun f2'(x, dx)
= let { y = x * x }
  ((x + y) * (x + x) + (x + y) * (x + x)) * dx
```

```
-----
Forward-mode derivative (CSE'd)
-----
```

```
fun f2'(x, dx)
= let { t1 = x + x * x }
  let { t2 = x + x }
  (t1 * t2 + t1 * t2) * dx
```

```
-----
Transposed Jacobian
-----
```

771	fun Rf2(x)	826
772	= lmTranspose(let { y = x * x }	827
773	lmScale((x + y) * (x + x) +	828
774	(x + y) * (x + x)))	829
775	-----	830
776	Optimised transposed Jacobian	831
777	-----	832
778	fun Rf2(x)	833
779	= let { y = x * x }	834
780	lmScale((x + y) * (x + x) +	835
781	(x + y) * (x + x))	836
782		837
783	-----	838
784	Reverse-mode derivative (unoptimised)	839
785	-----	840
786	fun f2'(x, dr)	841
787	= lmApply(let { y = x * x }	842
788	lmScale((x + y) * (x + x) +	843
789	(x + y) * (x + x)),	844
790	dr)	845
791	-----	846
792	Reverse-mode derivative (optimised)	847
793	-----	848
794	fun f2'(x, dr)	849
795	= let { y = x * x }	850
796	((x + y) * (x + x) +	851
797	(x + y) * (x + x)) * dr	852
798		853
799	-----	854
800	Reverse-mode derivative (CSE'd)	855
801	-----	856
802	fun f2'(x, dr)	857
803	= let { t1 = x + x * x }	858
804	let { t2 = x + x }	859
805	(t1 * t2 + t1 * t2) * dr	860
806		861
807		862
808		863
809		864
810		865
811		866
812		867
813		868
814		869
815		870
816		871
817		872
818		873
819		874
820		875
821		876
822		877
823		878
824		879
825		880