

# The simple essence of automatic differentiation

Conal Elliott

Target

January/June 2018

# Differentiable programming made easy

---

Current AI revolution runs on large data, speed, and AD, but

- AD algorithm (backprop) is complex and stateful.
- Graph APIs are complex and semantically dubious.

Solutions in this paper:

- AD: Simple, calculated, efficient, parallel-friendly, generalized.
- API: *derivative*.

# What's a derivative?

---

- Number
- Vector
- Covector
- Matrix
- Higher derivatives

Chain rule for each.

# What's a derivative?

$$\mathcal{D} :: (a \rightarrow b) \rightarrow (a \rightarrow (a \multimap b))$$

A local linear (affine) approximation:

$$\lim_{\varepsilon \rightarrow 0} \frac{\|f(a + \varepsilon) - (f a + \mathcal{D} f a \varepsilon)\|}{\|\varepsilon\|} = 0$$

See *Calculus on Manifolds* by Michael Spivak.

# Composition

Sequential:

$$\begin{aligned}(\circ) &:: (b \rightarrow c) \rightarrow (a \rightarrow b) \rightarrow (a \rightarrow c) \\(g \circ f) a &= g (f a)\end{aligned}$$

$$\mathcal{D} (g \circ f) a = \mathcal{D} g (f a) \circ \mathcal{D} f a \quad \text{-- chain rule}$$

Parallel:

$$\begin{aligned}(\triangle) &:: (a \rightarrow c) \rightarrow (a \rightarrow d) \rightarrow (a \rightarrow c \times d) \\(f \triangle g) a &= (f a, g a)\end{aligned}$$

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

# Linear functions

Linear functions are their own derivatives everywhere.

$$\mathcal{D} \text{ id } a = \text{id}$$

$$\mathcal{D} \text{ fst } a = \text{fst}$$

$$\mathcal{D} \text{ snd } a = \text{snd}$$

...

Chain rule:

$$\mathcal{D} (g \circ f) a = \mathcal{D} g (f a) \circ \mathcal{D} f a \quad \text{-- non-compositional}$$

To fix, combine regular result with derivative:

$$\hat{\mathcal{D}} :: (a \rightarrow b) \rightarrow (a \rightarrow (b \times (a \multimap b)))$$

$$\hat{\mathcal{D}} f = f \triangle \mathcal{D} f \quad \text{-- specification}$$

Often much work in common to  $f$  and  $\mathcal{D} f$ .

# Abstract algebra for functions

**class** *Category* ( $\rightsquigarrow$ ) **where**

*id* ::  $a \rightsquigarrow a$

( $\circ$ ) ::  $(b \rightsquigarrow c) \rightarrow (a \rightsquigarrow b) \rightarrow (a \rightsquigarrow c)$

**class** *Category* ( $\rightsquigarrow$ )  $\Rightarrow$  *Cartesian* ( $\rightsquigarrow$ ) **where**

*exl* ::  $(a \times b) \rightsquigarrow a$

*exr* ::  $(a \times b) \rightsquigarrow b$

( $\Delta$ ) ::  $(a \rightsquigarrow c) \rightarrow (a \rightsquigarrow d) \rightarrow (a \rightsquigarrow (c \times d))$

Plus laws and classes for arithmetic etc.



# Automatic differentiation

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

$\hat{D} :: (a \rightarrow b) \rightarrow D a b$

$\hat{D} f = D (f \triangle \mathcal{D} f)$  -- not computable

Specification:  $\hat{D}$  preserves *Category* and *Cartesian* structure:

$$\hat{D} id = id$$

$$\hat{D} (g \circ f) = \hat{D} g \circ \hat{D} f$$

$$\hat{D} exl = exl$$

$$\hat{D} exr = exr$$

$$\hat{D} (f \triangle g) = \hat{D} f \triangle \hat{D} g$$

*The game:* solve these equations for the RHS operations.

## Solution: simple automatic differentiation

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

*linearD*  $f = D\ (\lambda a \rightarrow (f\ a, f))$

**instance** *Category*  $D$  **where**

$id = linearD\ id$

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

**instance** *Cartesian*  $D$  **where**

$exl = linearD\ exl$

$exr = linearD\ exr$

$D\ f \triangle D\ g = D\ (\lambda a \rightarrow \mathbf{let}\ \{(b, f') = f\ a; (c, g') = g\ a\}\ \mathbf{in}\ ((b, c), f' \triangle g'))$

**instance** *NumCat*  $D$  **where**

$negate = linearD\ negate$

$add = linearD\ add$

$mul = D\ (mul \triangle (\lambda(a, b) \rightarrow \lambda(da, db) \rightarrow b * da + a * db))$

## Running examples

$sqr :: Num\ a \Rightarrow a \rightarrow a$

$sqr\ a = a * a$

$magSqr :: Num\ a \Rightarrow a \times a \rightarrow a$

$magSqr\ (a, b) = sqr\ a + sqr\ b$

$cosSinProd :: Floating\ a \Rightarrow a \times a \rightarrow a \times a$

$cosSinProd\ (x, y) = (cos\ z, sin\ z)$  **where**  $z = x * y$

In categorical vocabulary:

$sqr = mul \circ (id \triangle id)$

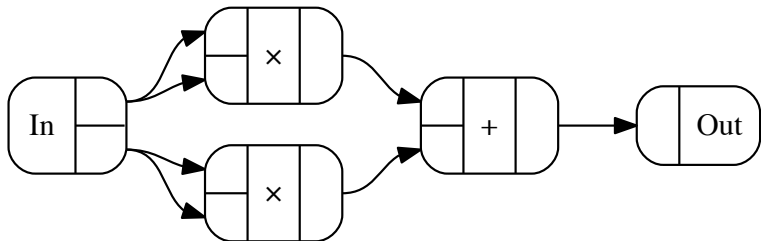
$magSqr = add \circ (mul \circ (exl \triangle exl) \triangle mul \circ (exr \triangle exr))$

$cosSinProd = (cos \triangle sin) \circ mul$

# Visualizing computations

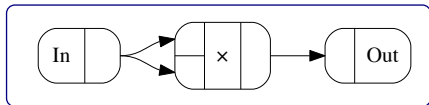
$$\text{magSqr } (a, b) = \text{sqr } a + \text{sqr } b$$

$$\text{magSqr} = \text{add} \circ (\text{mul} \circ (\text{exl} \triangle \text{exl}) \triangle \text{mul} \circ (\text{exr} \triangle \text{exr}))$$



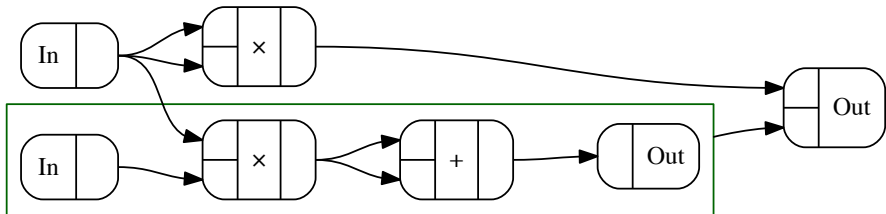
Auto-generated from Haskell code. See *Compiling to categories*.

# AD example



$$\text{sqr } a = a * a$$

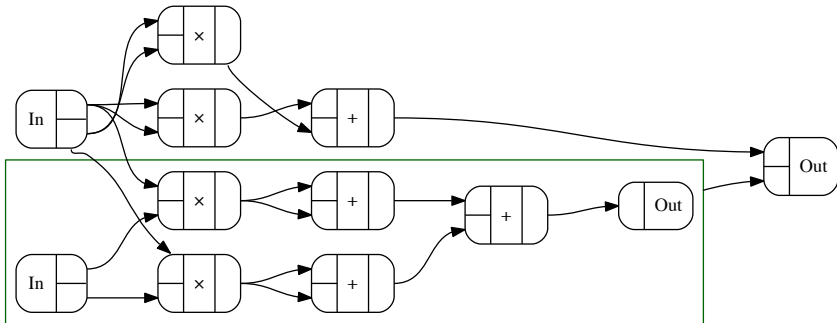
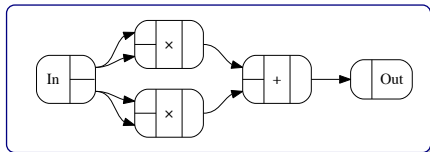
$$\text{sqr} = \text{mul} \circ (\text{id} \triangle \text{id})$$



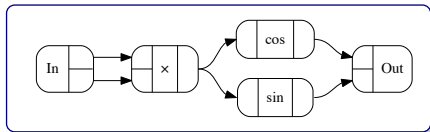
# AD example

$$\text{magSqr}(a, b) = \text{sqr } a + \text{sqr } b$$

$$\text{magSqr} = \text{add} \circ (\text{mul} \circ (\text{exl} \triangle \text{exl}) \triangle \text{mul} \circ (\text{exr} \triangle \text{exr}))$$

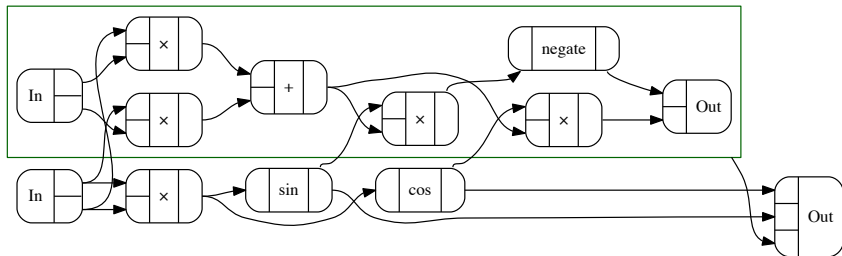


# AD example



$\text{cosSinProd}(x, y) = (\text{cos } z, \text{sin } z)$  **where**  $z = x * y$

$\text{cosSinProd} = (\text{cos } \Delta \text{ sin}) \circ \text{mul}$



# Generalizing AD

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

*linearD*  $f = D\ (\lambda a \rightarrow (f\ a, f))$

**instance** *Category*  $D$  **where**

$id = \text{linearD}\ id$

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

**instance** *Cartesian*  $D$  **where**

$exl = \text{linearD}\ exl$

$exr = \text{linearD}\ exr$

$D\ f \triangle D\ g = D\ (\lambda a \rightarrow \mathbf{let}\ \{(b, f') = f\ a; (c, g') = g\ a\}\ \mathbf{in}\ ((b, c), f' \triangle g'))$

Each  $D$  operation just uses corresponding  $(\multimap)$  operation.

Generalize from  $(\multimap)$  to other cartesian categories.



# Generalized AD

**newtype**  $D_{(\rightsquigarrow)}$   $a \ b = D (a \rightarrow b \times (a \rightsquigarrow b))$

*linearD*  $f \ f' = D (\lambda a \rightarrow (f \ a, f'))$

**instance** *Category*  $(\rightsquigarrow) \Rightarrow$  *Category*  $D_{(\rightsquigarrow)}$  **where**

$id = \text{linearD } id \ id$

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

**instance** *Cartesian*  $(\rightsquigarrow) \Rightarrow$  *Cartesian*  $D_{(\rightsquigarrow)}$  **where**

$exl = \text{linearD } exl \ exl$

$exr = \text{linearD } exr \ exr$

$D \ f \ \triangle \ D \ g = D (\lambda a \rightarrow \mathbf{let} \ \{(b, f') = f \ a; (c, g') = g \ a\} \ \mathbf{in} \ ((b, c), f' \ \triangle \ g'))$

**instance**  $\dots \Rightarrow$  *NumCat*  $D$  **where**

$negate = \text{linearD } negate \ negate$

$add = \text{linearD } add \ add$

$mul = ??$

# Numeric operations

Specific to (linear) *functions*:

$$mul = D (mul \triangle (\lambda(a, b) \rightarrow \lambda(da, db) \rightarrow b * da + a * db))$$

Rephrase:

$$scale :: \text{Multiplicative } a \Rightarrow a \rightarrow (a \circ a)$$

$$scale\ u = \lambda v \rightarrow u * v$$

$$(\nabla) :: (a \circ c) \rightarrow (b \circ c) \rightarrow ((a \times b) \circ c)$$

$$f \nabla g = \lambda(a, b) \rightarrow f\ a + g\ b$$

Now

$$mul = D (mul \triangle (\lambda(a, b) \rightarrow scale\ b \nabla scale\ a))$$

# Linear arrow (biproduct) vocabulary

**class** *Category* ( $\rightsquigarrow$ ) **where**

*id* ::  $a \rightsquigarrow a$

( $\circ$ ) ::  $(b \rightsquigarrow c) \rightarrow (a \rightsquigarrow b) \rightarrow (a \rightsquigarrow c)$

**class** *Category* ( $\rightsquigarrow$ )  $\Rightarrow$  *Cartesian* ( $\rightsquigarrow$ ) **where**

*exl* ::  $(a \times b) \rightsquigarrow a$

*exr* ::  $(a \times b) \rightsquigarrow b$

( $\triangle$ ) ::  $(a \rightsquigarrow c) \rightarrow (a \rightsquigarrow d) \rightarrow (a \rightsquigarrow (c \times d))$

**class** *Category* ( $\rightsquigarrow$ )  $\Rightarrow$  *Cocartesian* ( $\rightsquigarrow$ ) **where**

*inl* ::  $a \rightsquigarrow (a \times b)$

*inr* ::  $b \rightsquigarrow (a \times b)$

( $\nabla$ ) ::  $(a \rightsquigarrow c) \rightarrow (b \rightsquigarrow c) \rightarrow ((a \times b) \rightsquigarrow c)$

**class** *ScalarCat* ( $\rightsquigarrow$ ) *a* **where**

*scale* ::  $a \rightarrow (a \rightsquigarrow a)$

# Linear transformations as functions

**newtype**  $a \rightarrow^+ b = \text{AddFun } (a \rightarrow b)$

**instance** *Category*  $(\rightarrow^+)$  **where**

$id = \text{AddFun } id$

$(\circ) = \text{inNew}_2 (\circ)$

**instance** *Cartesian*  $(\rightarrow^+)$  **where**

$exl = \text{AddFun } exl$

$exr = \text{AddFun } exr$

$(\Delta) = \text{inNew}_2 (\Delta)$

**instance** *Cocartesian*  $(\rightarrow^+)$  **where**

$inl = \text{AddFun } (, 0)$

$inr = \text{AddFun } (0, )$

$(\nabla) = \text{inNew}_2 (\lambda f g (x, y) \rightarrow f x + g y)$

**instance** *Multiplicative*  $s \Rightarrow \text{ScalarCat } (\rightarrow^+) s$  **where**

$scale s = \text{AddFun } (s *)$

# Extracting a data representation

- How to extract a matrix or gradient vector?
- Sample over a domain *basis* (rows of identity matrix).
- For  $n$ -dimensional *domain*,
  - Make  $n$  passes.
  - Each pass works on  $n$ -D sparse (“one-hot”) input.
  - Very inefficient.
- For gradient-based optimization,
  - High-dimensional domain.
  - Very low-dimensional (1-D) codomain.

# Generalized matrices

**newtype**  $M_s$   $a$   $b = L (V_s b (V_s a s))$

*applyL* ::  $M_s$   $a$   $b \rightarrow (a \rightarrow b)$

Require *applyL* to preserve structure. Solve for methods.

# Core vocabulary

Sufficient to build arbitrary “matrices”:

$scale :: a \rightarrow (a \rightsquigarrow a)$  --  $1 \times 1$

$(\nabla) :: (a \rightsquigarrow c) \rightarrow (b \rightsquigarrow c) \rightarrow ((a \times b) \rightsquigarrow c)$  -- horizontal juxt

$(\Delta) :: (a \rightsquigarrow c) \rightarrow (a \rightsquigarrow d) \rightarrow (a \rightsquigarrow (c \times d))$  -- vertical juxt

Types guarantee rectangularity.

# Efficiency of composition

- Arrow composition is associative.
- Some associations are more efficient than others, so
  - Associate optimally.
  - Equivalent to *matrix chain multiplication* —  $O(n \log n)$ .
  - Choice determined by *types*, i.e., compile-time information.
- All-right: “forward mode AD” (FAD).
- All-left: “reverse mode AD” (RAD).
- RAD is *much* better for gradient-based optimization.



# Left-associating composition (RAD)

- CPS-like category:
  - Represent  $a \rightsquigarrow b$  by  $(b \rightsquigarrow r) \rightarrow (a \rightsquigarrow r)$ .
  - Meaning:  $f \mapsto (\circ f)$ .
  - Results in left-composition.
  - Initialize with  $id :: r \rightsquigarrow r$ .
  - Construct  $h \circ \mathcal{D} f a$  directly, without  $\mathcal{D} f a$ .

# Continuation category

**newtype**  $Cont_{(\rightsquigarrow)}^r a b = Cont ((b \rightsquigarrow r) \rightarrow (a \rightsquigarrow r))$

$cont :: Category (\rightsquigarrow) \Rightarrow (a \rightsquigarrow b) \rightarrow Cont_{(\rightsquigarrow)}^r a b$

$cont f = Cont (\circ f)$

Require  $cont$  to preserve structure. Solve for methods.

We'll use an isomorphism:

$join :: Cocartesian (\rightsquigarrow) \Rightarrow (c \rightsquigarrow a) \times (d \rightsquigarrow a) \rightarrow ((c \times d) \rightsquigarrow a)$

$unjoin :: Cocartesian (\rightsquigarrow) \Rightarrow ((c \times d) \rightsquigarrow a) \rightarrow (c \rightsquigarrow a) \times (d \rightsquigarrow a)$

$join (f, g) = f \nabla g$

$unjoin h = (h \circ inl, h \circ inr)$

# Continuation category (solution)

**instance**  $\text{Category } (\rightsquigarrow) \Rightarrow \text{Category } \text{Cont}_{(\rightsquigarrow)}^r$  **where**

$\text{id} = \text{Cont id}$

$\text{Cont } g \circ \text{Cont } f = \text{Cont } (f \circ g)$

**instance**  $\text{Cartesian } (\rightsquigarrow) \Rightarrow \text{Cartesian } \text{Cont}_{(\rightsquigarrow)}^r$  **where**

$\text{exl} = \text{Cont } (\text{join} \circ \text{inl})$

$\text{exr} = \text{Cont } (\text{join} \circ \text{inr})$

$(\triangle) = \text{inNew}_2 (\lambda f \ g \rightarrow (f \nabla g) \circ \text{unjoin})$

**instance**  $\text{Cocartesian } (\rightsquigarrow) \Rightarrow \text{Cocartesian } \text{Cont}_{(\rightsquigarrow)}^r$  **where**

$\text{inl} = \text{Cont } (\text{exl} \circ \text{unjoin})$

$\text{inr} = \text{Cont } (\text{exr} \circ \text{unjoin})$

$(\nabla) = \text{inNew}_2 (\lambda f \ g \rightarrow \text{join} \circ (f \triangle g))$

**instance**  $\text{ScalarCat } (\rightsquigarrow) \ a \Rightarrow \text{ScalarCat } \text{Cont}_{(\rightsquigarrow)}^r \ a$  **where**

$\text{scale } s = \text{Cont } (\text{scale } s)$

# Reverse-mode AD without tears

---

$$DCont_{M_S}^r$$

# Left-associating composition (RAD)

- CPS-like category:
  - Represent  $a \rightsquigarrow b$  by  $(b \rightsquigarrow r) \rightarrow (a \rightsquigarrow r)$ .
  - Meaning:  $f \mapsto (\circ f)$ .
  - Results in left-composition.
  - Initialize with  $id :: r \rightsquigarrow r$ .
  - Construct  $h \circ \mathcal{D} f a$  directly, without  $\mathcal{D} f a$ .
- We've seen this trick before:
  - Transforming naive *reverse* from quadratic to linear.
  - List generalizes to monoids, and monoids to categories.

## One of my favorite papers

- *Continuation-Based Program Transformation Strategies*  
Mitch Wand, 1980, JACM.
- Introduce a continuation argument, e.g.,  $[a] \rightarrow [a]$ .
- Notice the continuations that arise, e.g.,  $(\# as)$ .
- Find a *data* representation, e.g.,  $as :: [a]$
- Identify associative operation that represents composition, e.g.,  $(\#)$ , since  $(\# bs) \circ (\# as) = (\# (as \# bs))$ .

- Vector space dual:  $u \multimap s$ , with  $u$  a vector space over  $s$ .
- If  $u$  has finite dimension, then  $u \multimap s \cong u$ .
- For  $f :: u \multimap s$ ,  $f = \text{dot } v$  for some  $v :: u$ .
- Gradients are derivatives of functions with scalar codomain.
- Represent  $a \multimap b$  by  $(b \multimap s) \rightarrow (a \multimap s)$  by  $b \rightarrow a$ .
- *Ideal* for extracting gradient vector. Just apply to 1 (*id*).

**newtype**  $Dual_{(\rightsquigarrow)} a b = Dual (b \rightsquigarrow a)$

$asDual :: Cont_{(\rightsquigarrow)}^s a b \rightarrow Dual_{(\rightsquigarrow)} a b$

$asDual (Cont f) = Dual (dot^{-1} \circ f \circ dot)$

where

$dot :: u \rightarrow (u \multimap s)$

$dot^{-1} :: (u \multimap s) \rightarrow u$

Require  $asDual$  to preserve structure. Solve for methods.



## Duality (solution)

**newtype**  $Dual_{(\rightsquigarrow)}$   $a\ b = Dual\ (b \rightsquigarrow a)$

**instance**  $Category\ (\rightsquigarrow) \Rightarrow Category\ Dual_{(\rightsquigarrow)}$  **where**

$id = Dual\ id$

$(\circ) = inNew_2\ (flip\ \circ)$

**instance**  $Cocartesian\ (\rightsquigarrow) \Rightarrow Cartesian\ Dual_{(\rightsquigarrow)}$  **where**

$exl = Dual\ inl$

$exr = Dual\ inr$

$(\Delta) = inNew_2\ (\nabla)$

**instance**  $Cartesian\ (\rightsquigarrow) \Rightarrow Cocartesian\ Dual_{(\rightsquigarrow)}$  **where**

$inl = Dual\ exl$

$inr = Dual\ exr$

$(\nabla) = inNew_2\ (\Delta)$

**instance**  $ScalarCat\ (\rightsquigarrow)\ s \Rightarrow ScalarCat\ Dual_{(\rightsquigarrow)}\ s$  **where**

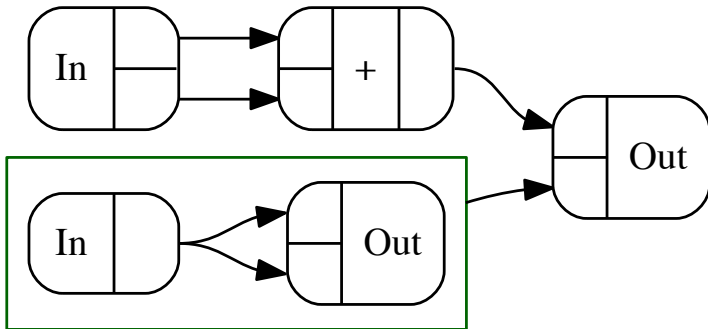
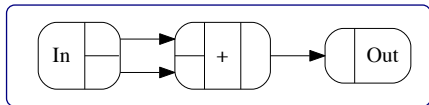
$scale\ s = Dual\ (scale\ s)$

# Backpropagation

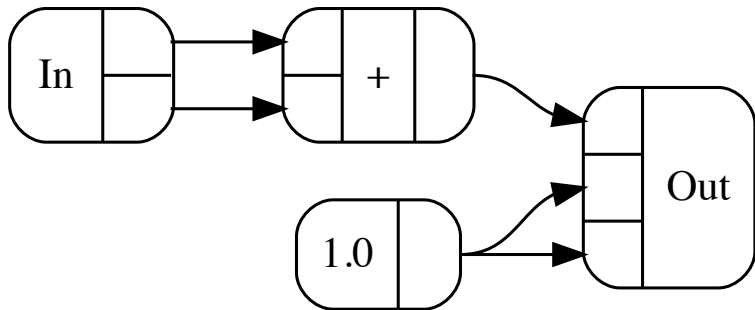
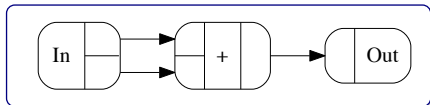
---

$$D_{Dual_{\rightarrow+}}$$

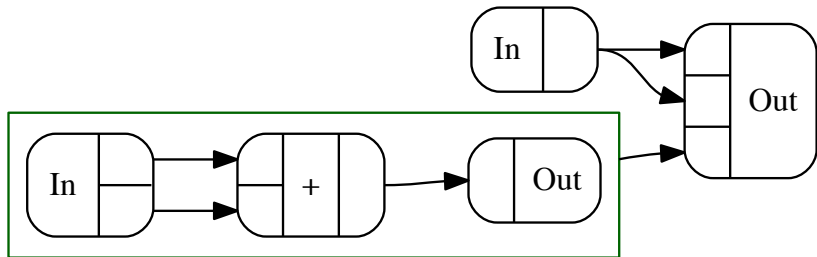
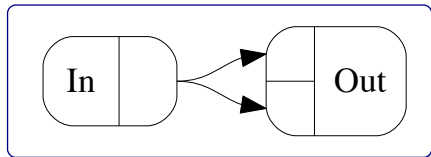
# RAD example (dual function)



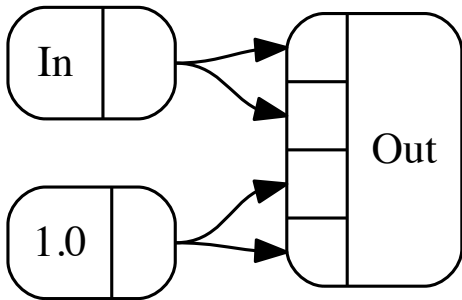
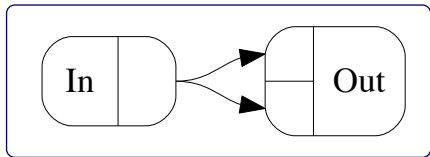
# RAD example (dual vector)



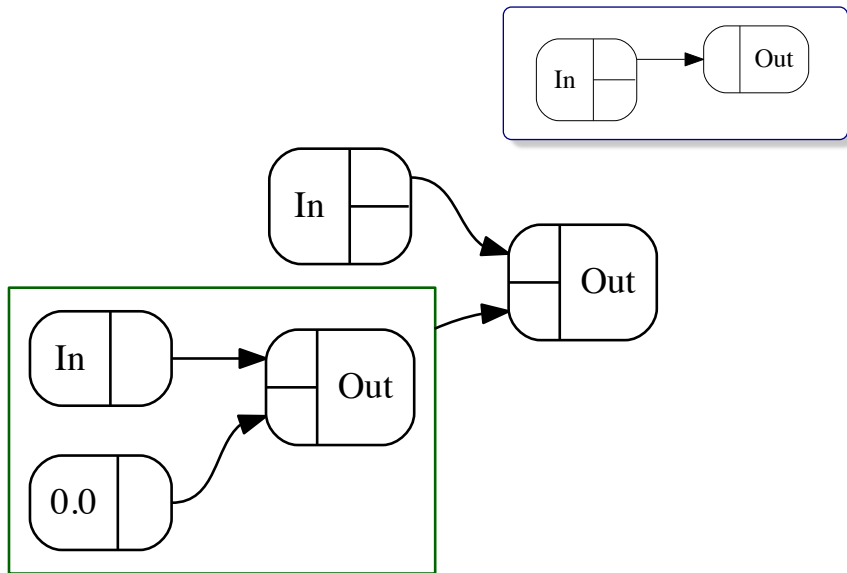
## RAD example (dual function)



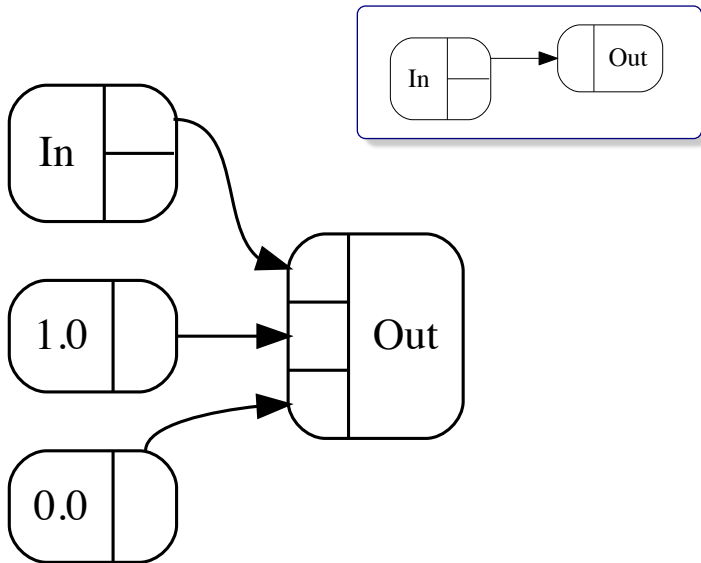
# RAD example (vector)



# RAD example (dual function)

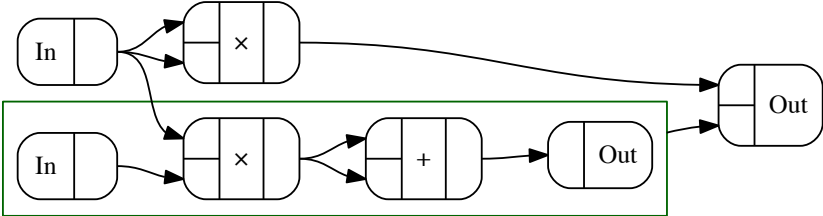
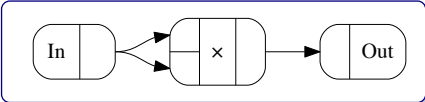


# RAD example (dual vector)

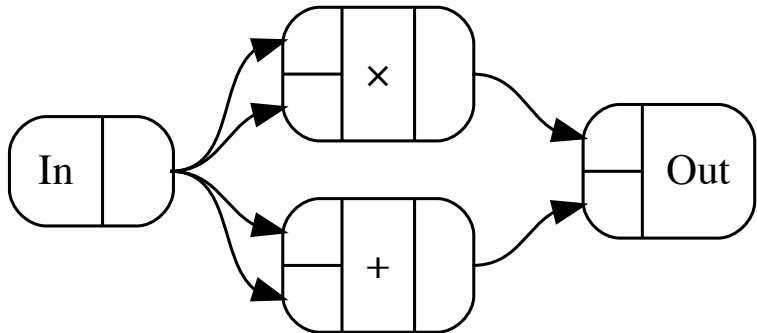
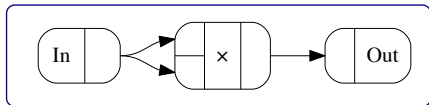




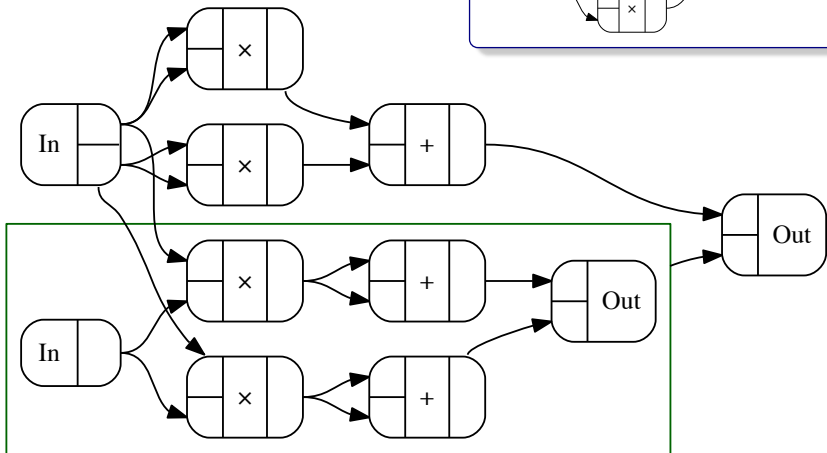
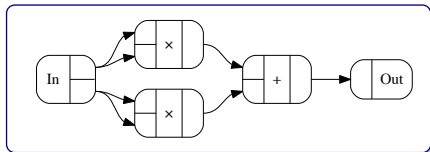
# RAD example (dual function)



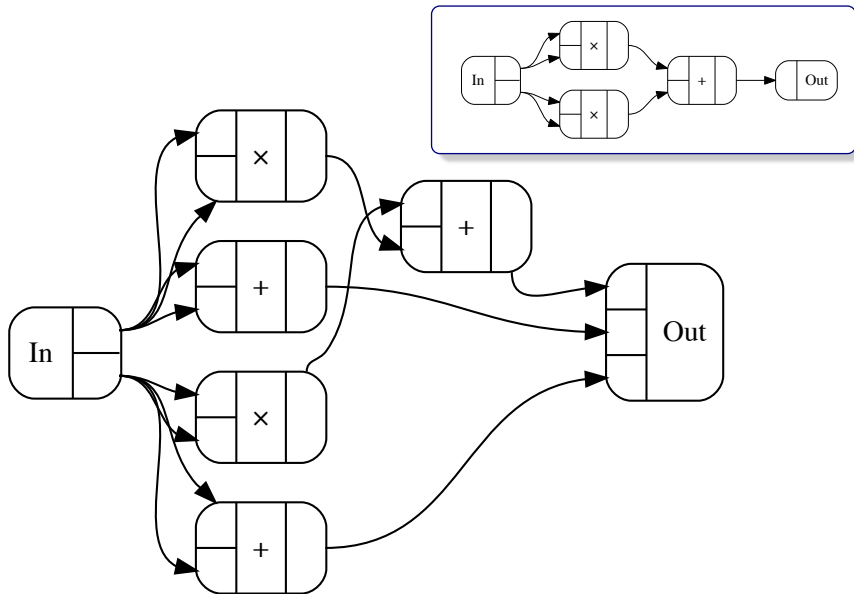
# RAD example (dual vector)



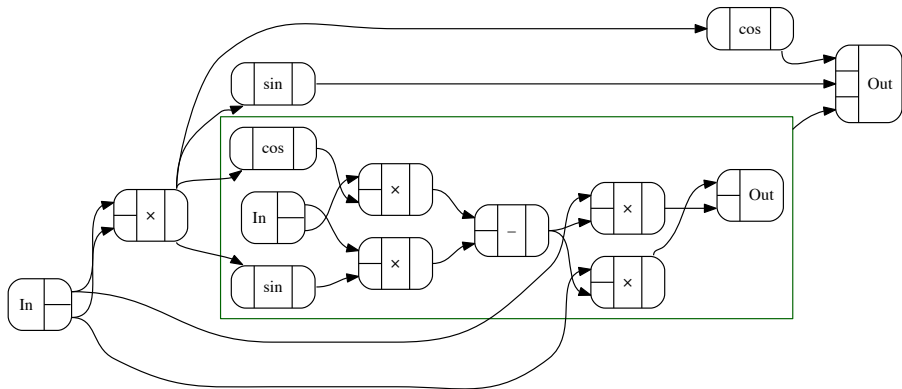
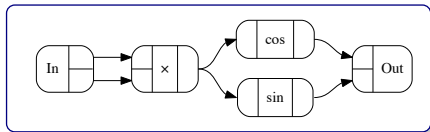
# RAD example (dual function)



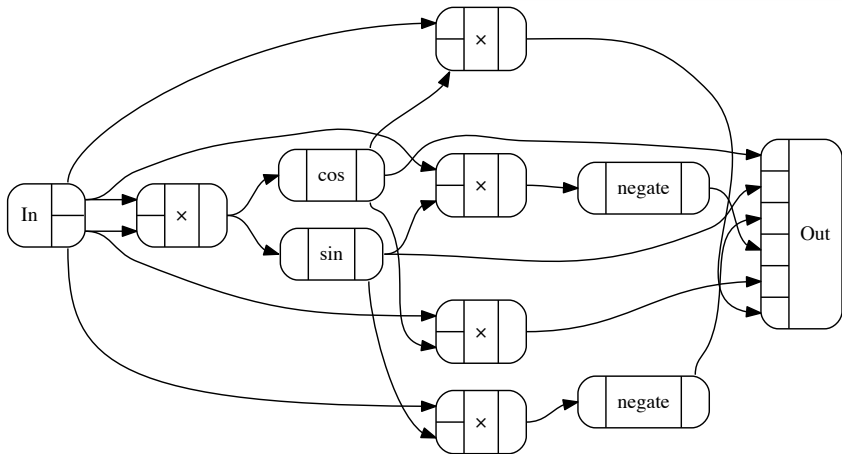
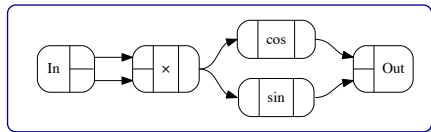
# RAD example (dual vector)



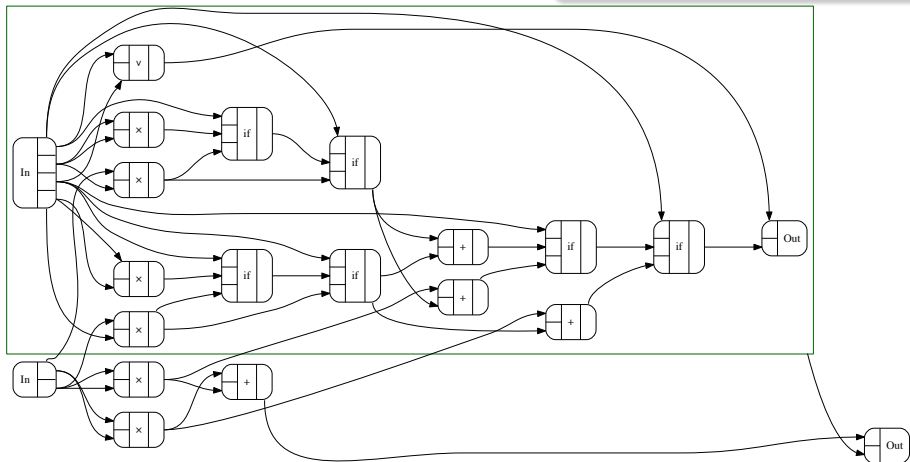
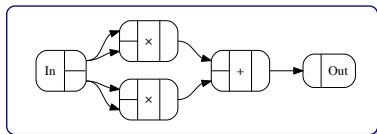
# RAD example (dual function)



# RAD example (matrix)



# Incremental evaluation



# Conclusions

---

- Simple AD algorithm, specializing to forward, reverse, mixed.
- No graphs, tapes, tags, partial derivatives, or mutation.
- Parallel-friendly and low memory use.
- Calculated from simple, regular algebra problems.
- Generalizes to derivative categories other than linear maps.
- Differentiate regular Haskell code (via plugin).
- More details in my [ICFP 2018 paper](#).



## Reflections: recipe for success

---

Key principles:

- Capture main concepts as first-class values.
- Focus on abstract notions, not specific representations.
- Calculate efficient implementation from simple specification.

Not previously applied to AD (afaik).

*Quandary:* Most programming languages poor for function-like things.

*Solution: Compiling to categories.*

# Symbolic vs automatic differentiation

Often described as opposing techniques:

- *Symbolic*:
  - Apply differentiation rules symbolically.
  - Can duplicate much work.
  - Needs algebraic manipulation.
- *Automatic*:
  - FAD: easy to implement but often inefficient.
  - RAD: efficient but tricky to implement.

My view: *AD is SD done by a compiler.*

Compilers already work symbolically and preserve sharing.