

# Microsoft Research

Each year Microsoft Research hosts hundreds of influential speakers from around the world including leading scientists, renowned experts in technology, book authors, and leading academics, and makes videos of these lectures freely available.

2016 © Microsoft Corporation. All rights reserved.

# The simple essence of automatic differentiation

Conal Elliott

Target

January/June 2018

# What's a derivative?

- Number
- Vector
- Covector
- Matrix
- Higher derivatives

Chain rule for each.

# What's a derivative?

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

where

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

See *Calculus on Manifolds* by Michael Spivak.

# Composition

Sequential:

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

$$(g \circ f) a = g (f a)$$

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

# What's a derivative?

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

where

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

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:

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

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



# Compositionality

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:

$$\begin{aligned} \hat{\mathcal{D}} &:: (a \rightarrow b) \rightarrow (a \rightarrow (b \times (a \multimap b))) \\ \hat{\mathcal{D}} f &= f \triangle \mathcal{D} f \quad \text{-- specification} \end{aligned}$$

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



# Linear functions

Linear functions are their own perfect linear approximations.

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

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

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

...

For linear functions  $f$ ,

$$\hat{\mathcal{D}} f a = (f a, f)$$

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

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

✦

Plus laws and classes for arithmetic etc.

# Linear functions

Linear functions are their own perfect linear approximations.

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

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

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

...

For linear functions  $f$ ,

$$\hat{\mathcal{D}} f a = (f a, f)$$

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

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

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

Require  $\hat{D}$  to preserve *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**  $D\ a\ b = D\ (a \rightarrow b \times (a \multimap b))$

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

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

$id = linearD\ id$

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

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

$exl = linearD\ exl$

$exr = linearD\ exr$

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

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

$negate = linearD\ negate$

$add = linearD\ add$

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

# 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

Require  $\hat{D}$  to preserve *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**  $D\ a\ b = D\ (a \rightarrow b \times (a \multimap b))$

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

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

$id = linearD\ id$

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

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

$exl = linearD\ exl$

$exr = linearD\ exr$

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

**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

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

$sqr\ a = a * a$

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

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

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

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

In categorical vocabulary:

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

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

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

## Solution: simple automatic differentiation

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

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

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

$id = linearD\ id$

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

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

$exl = linearD\ exl$

$exr = linearD\ exr$

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

**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

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

$sqr\ a = a * a$

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

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

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

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

In categorical vocabulary:

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

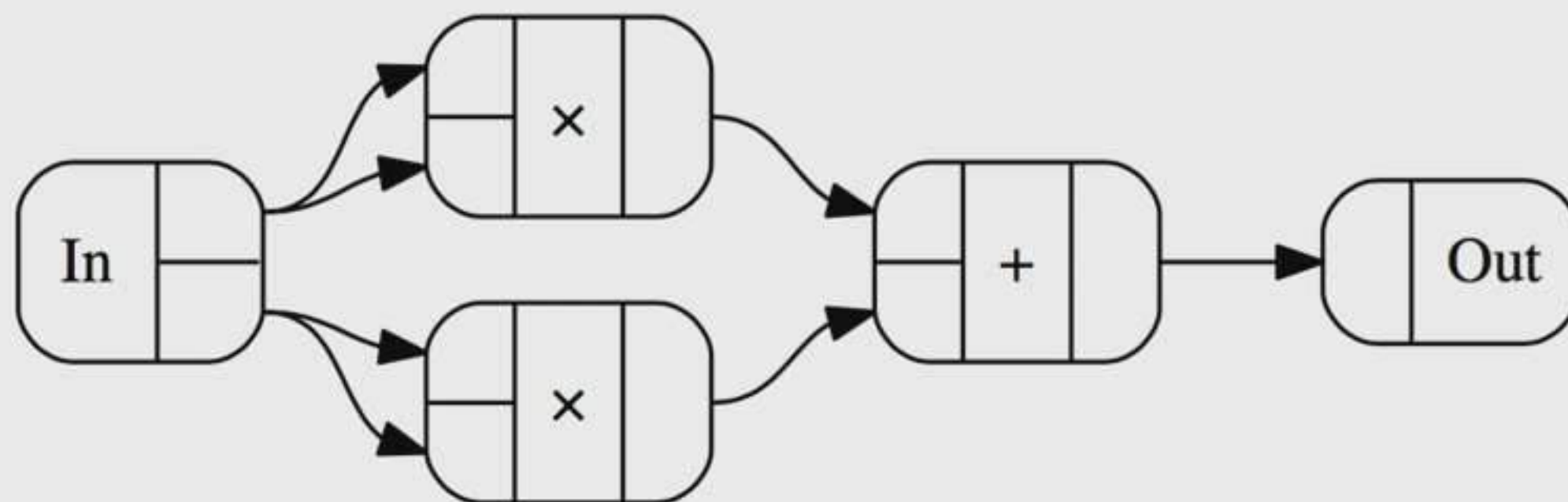
$magSqr = add \circ (mul \circ (exl \triangle exl) \triangle mul \circ (exr \triangle exr))$

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

# Visualizing computations

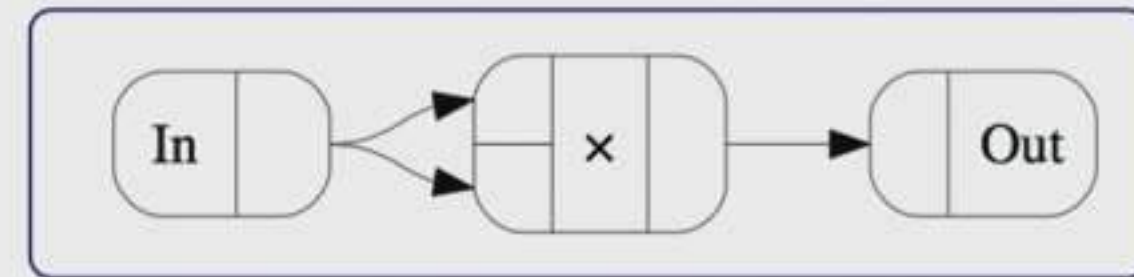
$$\text{magSqr } (a, b) = \text{sqr } a + \text{sqr } b$$

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



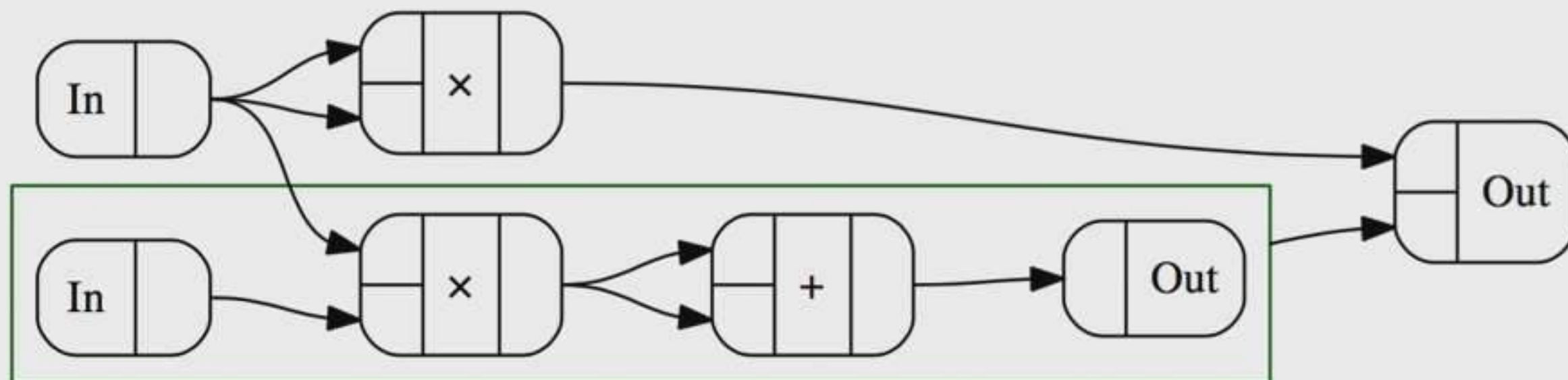
Auto-generated from Haskell code. See *Compiling to categories*.

# AD example



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

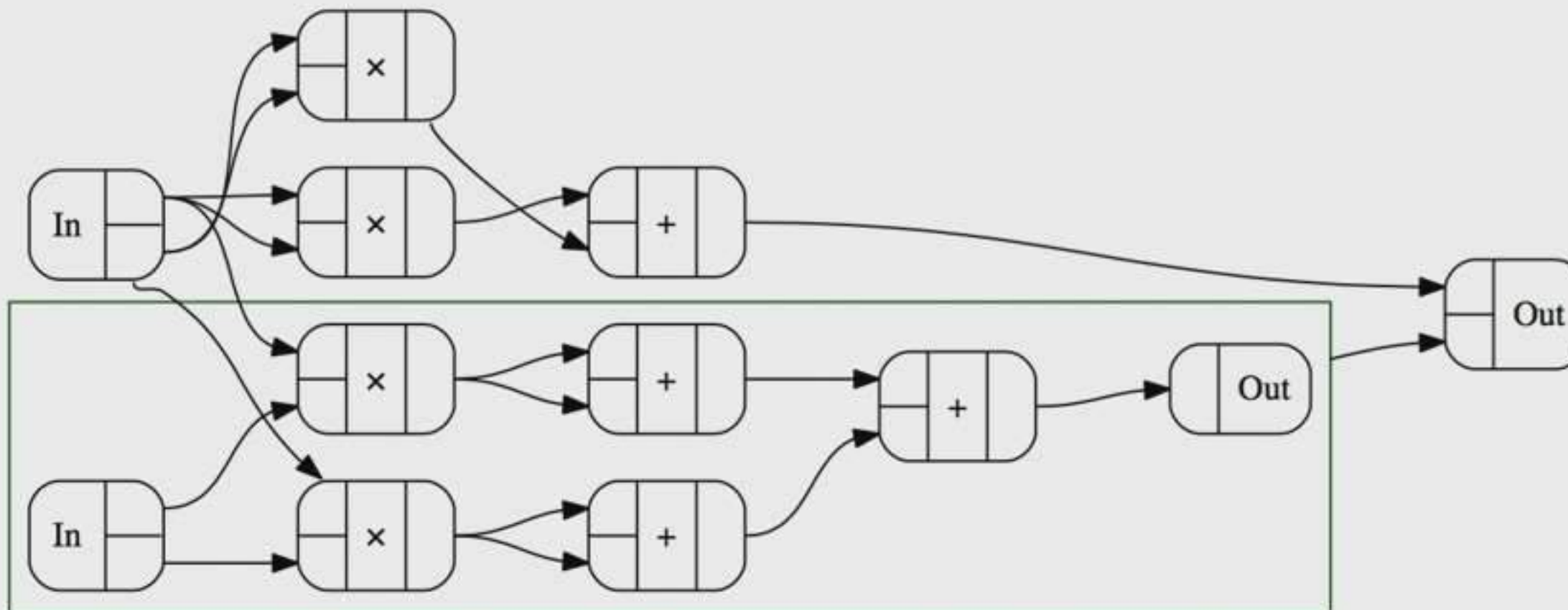
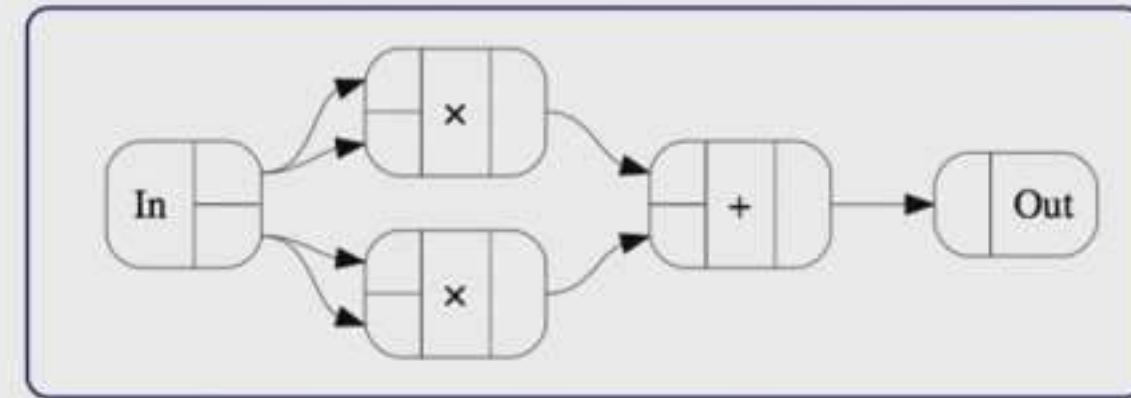
$$\text{sqr} = \text{mul} \circ (\text{id} \triangle \text{id})$$



# AD example

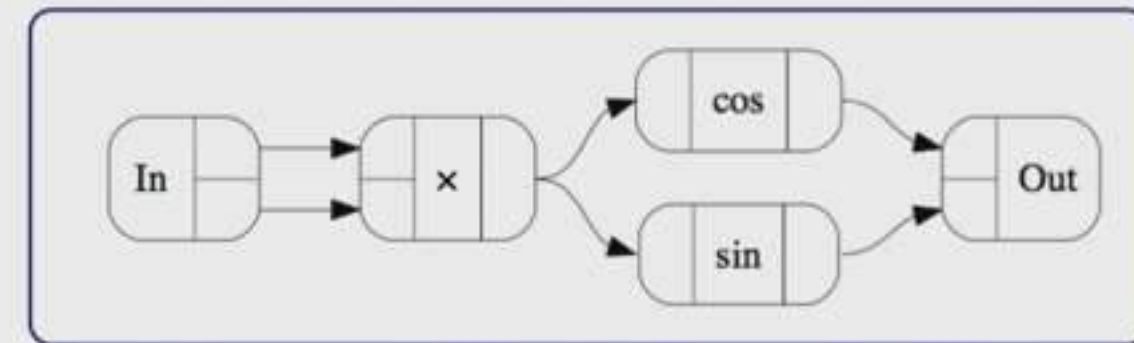
$$\text{magSqr}(a, b) = \text{sqr } a + \text{sqr } b$$

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



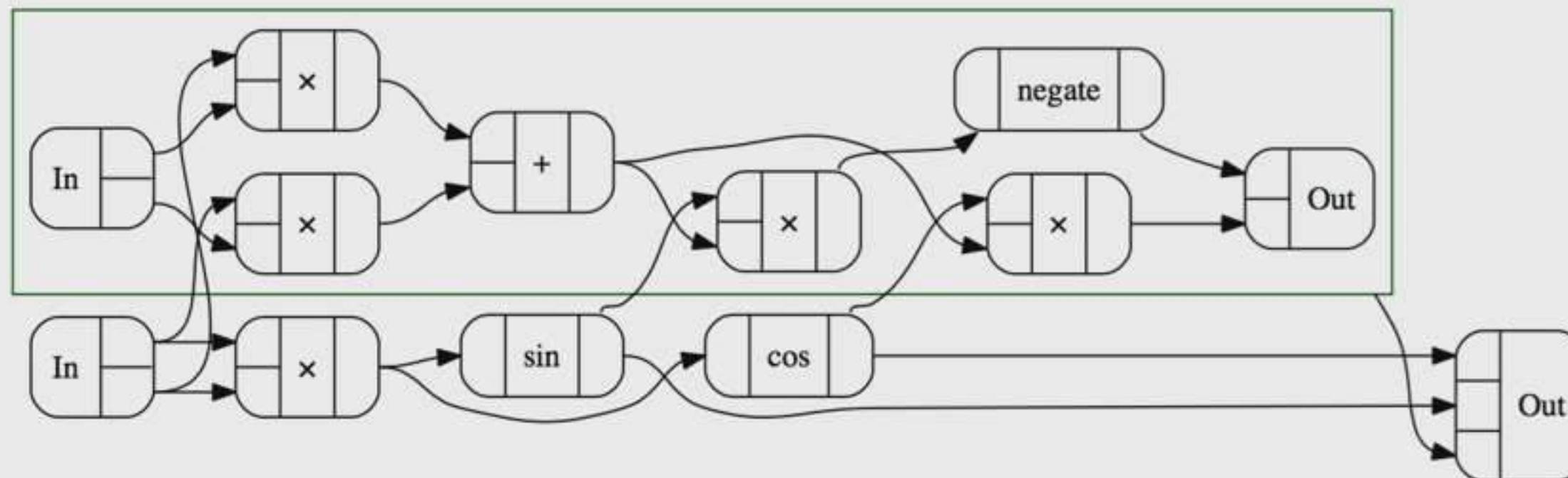


# AD example



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

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



# 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**

$id = linearD\ id$

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

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

$exl = linearD\ exl$

$exr = linearD\ exr$

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

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

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

# Generalized AD

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

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

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

$id = linearD id id$

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

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

$exl = linearD exl exl$

$exr = linearD exr exr$

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

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

$negate = linearD negate negate$

$add = linearD add add$

$mul = ??$



# 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**

$id = linearD\ id$

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

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

$exl = linearD\ exl$

$exr = linearD\ exr$

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

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* = *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* = *linearD exl exl*

*exr* = *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**...  $\Rightarrow NumCat\ D$  **where**

*negate* = *linearD negate negate*

*add* = *linearD add add*

*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**

$id = linearD\ id$

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

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

$exl = linearD\ exl$

$exr = linearD\ exr$

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

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

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



# Generalized AD

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

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

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

$id = linearD \ id \ id$

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

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

$exl = linearD \ exl \ exl$

$exr = linearD \ exr \ exr$

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

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

$negate = linearD \ negate \ negate$

$add = linearD \ add \ add$

$mul = ??$



# Numeric operations

Specific to (linear) *functions*:

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

Rephrase:

$$scale :: Multiplicative\ a \Rightarrow a \rightarrow (a \multimap a)$$

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

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

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

Now

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

# Linear arrow vocabulary

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

$id :: a \rightsquigarrow a$

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

**class** *Category*  $(\rightsquigarrow) \Rightarrow \text{Cartesian } (\rightsquigarrow)$  **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*  $(\rightsquigarrow) \Rightarrow \text{Cocartesian } (\rightsquigarrow)$  **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*  $(\rightsquigarrow) \text{ } a$  **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**

$exl = \text{AddFun } exl$

$exr = \text{AddFun } exr$

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

**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”:

$$scale :: a \rightarrow (a \rightsquigarrow a) \quad \text{-- } 1 \times 1$$

$$(\nabla) \quad :: (a \rightsquigarrow c) \rightarrow (b \rightsquigarrow c) \rightarrow ((a \times b) \rightsquigarrow c) \quad \text{-- horizontal juxt}$$

$$(\Delta) \quad :: (a \rightsquigarrow c) \rightarrow (a \rightsquigarrow d) \rightarrow (a \rightsquigarrow (c \times d)) \quad \text{-- vertical juxt}$$

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**  $Category (\rightsquigarrow) \Rightarrow Category\ Cont^r_{(\rightsquigarrow)}$  **where**

$id = Cont\ id$

$Cont\ g \circ Cont\ f = Cont\ (f \circ g)$

**instance**  $Cartesian (\rightsquigarrow) \Rightarrow Cartesian\ Cont^r_{(\rightsquigarrow)}$  **where**

$exl = Cont\ (join \circ inl)$

$exr = Cont\ (join \circ inr)$

$(\triangle) = inNew_2\ (\lambda f\ g \rightarrow (f \nabla g) \circ unjoin)$

**instance**  $Cocartesian (\rightsquigarrow) \Rightarrow Cocartesian\ Cont^r_{(\rightsquigarrow)}$  **where**

$inl = Cont\ (exl \circ unjoin)$

$inr = Cont\ (exr \circ unjoin)$

$(\nabla) = inNew_2\ (\lambda f\ g \rightarrow join \circ (f \triangle g))$

**instance**  $ScalarCat (\rightsquigarrow) a \Rightarrow ScalarCat\ Cont^r_{(\rightsquigarrow)} a$  **where**

$scale\ s = Cont\ (scale\ s)$

# Reverse-mode AD without tears

$$D_{Cont}^r_{M_s}$$



# Duality

- 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*).

# Duality

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

$asDual :: Cont^s_{(\rightsquigarrow)} a\ b \rightarrow Dual_{(\rightsquigarrow)} a\ b$

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

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**

$exl = Dual\ inl$

$exr = Dual\ inr$

$(\triangle) = inNew_2\ (\nabla)$

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

$inl = Dual\ exl$

$inr = Dual\ exr$

$(\nabla) = inNew_2\ (\triangle)$

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

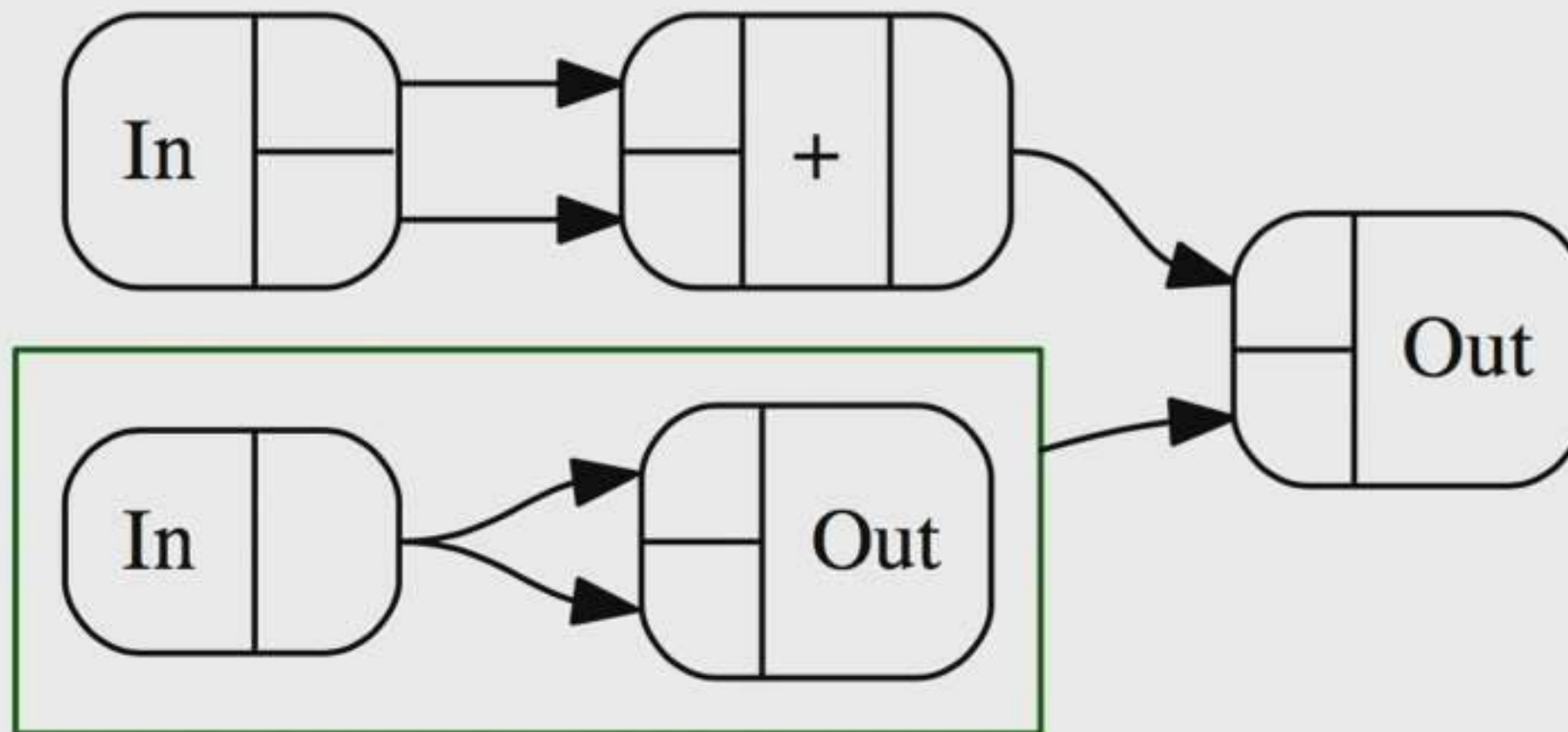
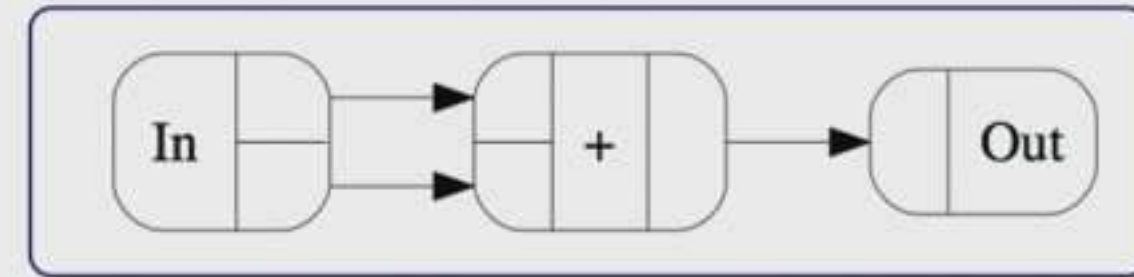
$scale\ s = Dual\ (scale\ s)$

# Backpropagation

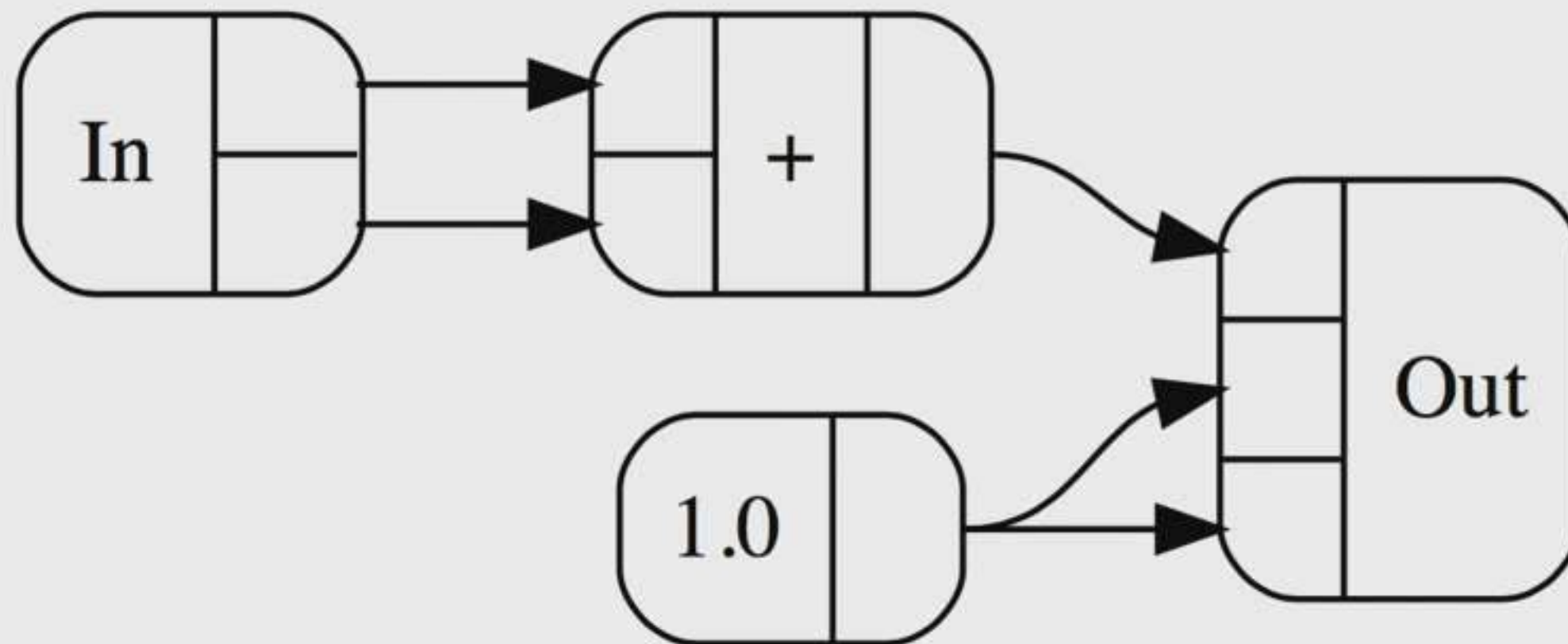
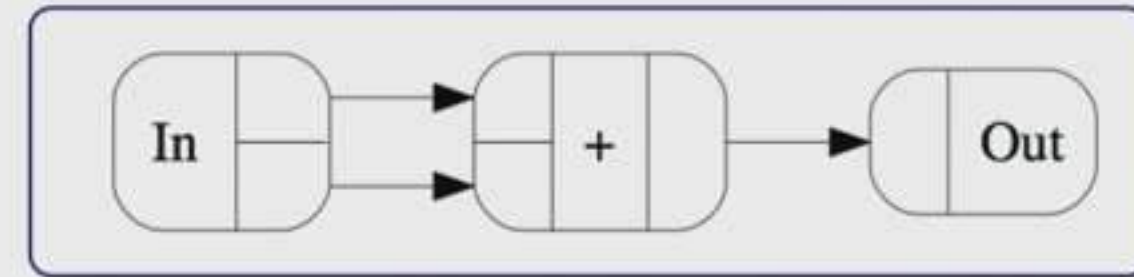
$$D_{Dual_{\rightarrow+}}$$



# RAD example (dual function)



## RAD example (dual vector)



## Solution: simple automatic differentiation

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

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

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

$id = linearD\ id$

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

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

$exl = linearD\ exl$

$exr = linearD\ exr$

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

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

$negate = linearD\ negate$

$add = linearD\ add$

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

# 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**

$id = linearD\ id$

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

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

$exl = linearD\ exl$

$exr = linearD\ exr$

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

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

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