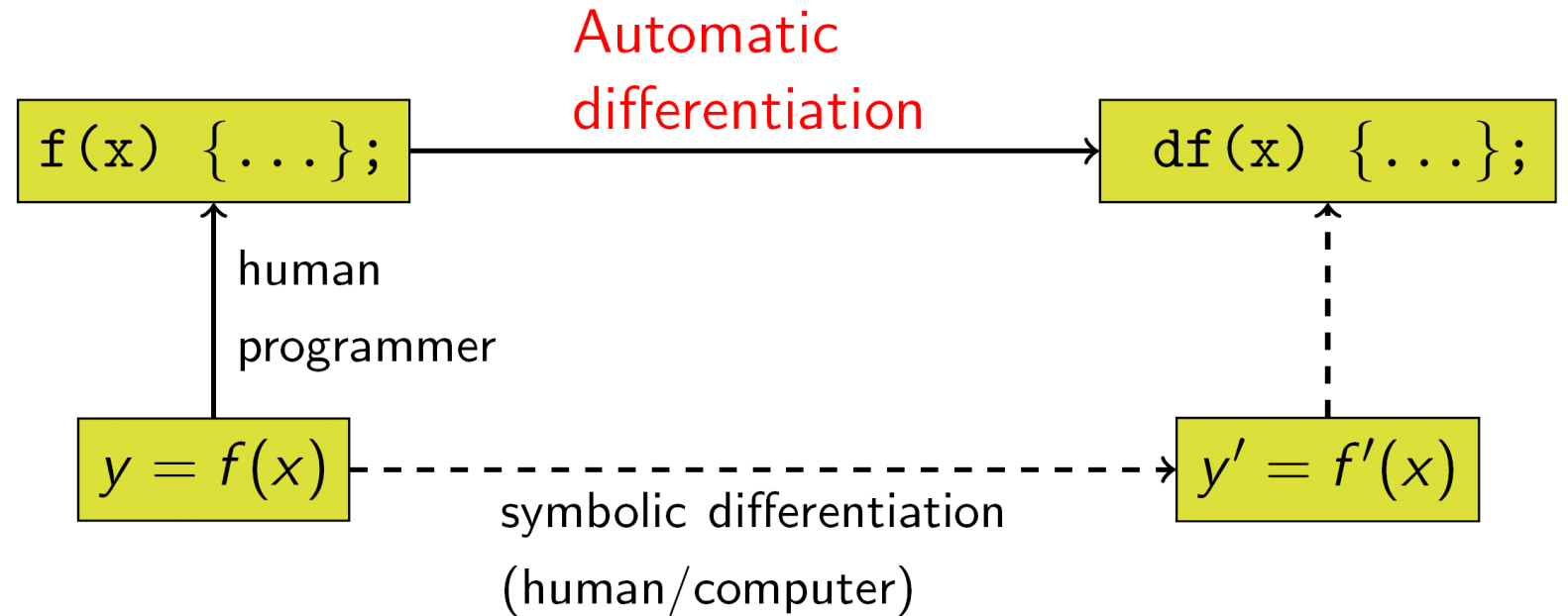


Efficient dual numbers reverse automatic differentiation by well-known program transformations

Paper by Tom Smeding and Matthijs Vákár
Presentation by Simon van Hus

Automatic Differentiation

- Manual Differentiation
- Symbolic Differentiation
- Numeric Differentiation
- Automatic Differentiation



Automatic Differentiation

- Manual Differentiation
- Symbolic Differentiation
- Numeric Differentiation
- Automatic Differentiation
 - **Forward Mode**

$$\begin{aligned}\frac{\partial y}{\partial x} &= \frac{\partial y}{\partial w_{n-1}} \frac{\partial w_{n-1}}{\partial x} \\ &= \frac{\partial y}{\partial w_{n-1}} \left(\frac{\partial w_{n-1}}{\partial w_{n-2}} \frac{\partial w_{n-2}}{\partial x} \right) \\ &= \frac{\partial y}{\partial w_{n-1}} \left(\frac{\partial w_{n-1}}{\partial w_{n-2}} \left(\frac{\partial w_{n-2}}{\partial w_{n-3}} \frac{\partial w_{n-3}}{\partial x} \right) \right) \\ &= \dots\end{aligned}$$

Automatic Differentiation

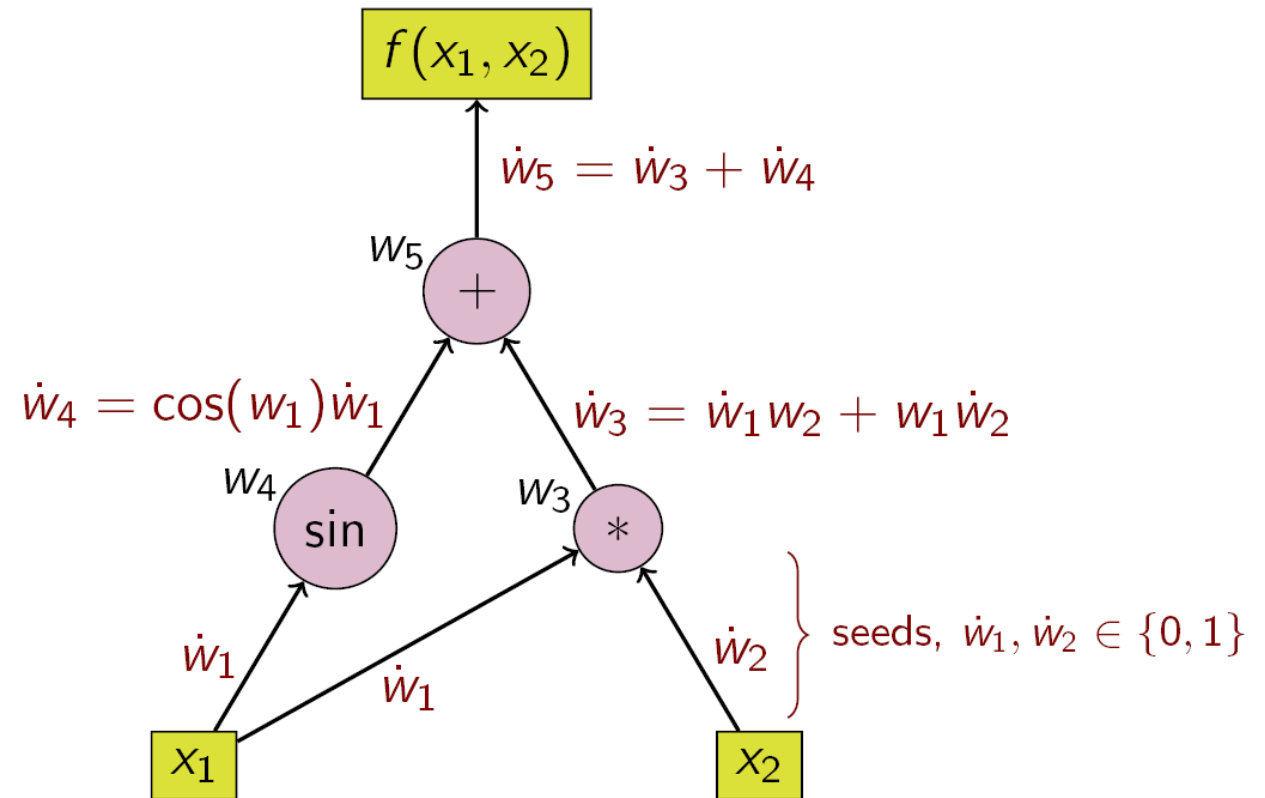
- Manual Differentiation
- Symbolic Differentiation
- Numeric Differentiation
- Automatic Differentiation
 - **Forward Mode**
 - **Reverse Mode**

$$\begin{aligned}\frac{\partial y}{\partial x} &= \frac{\partial y}{\partial w_1} \frac{\partial w_1}{\partial x} \\ &= \left(\frac{\partial y}{\partial w_2} \frac{\partial w_2}{\partial w_1} \right) \frac{\partial w_1}{\partial x} \\ &= \left(\left(\frac{\partial y}{\partial w_3} \frac{\partial w_3}{\partial w_2} \right) \frac{\partial w_2}{\partial w_1} \right) \frac{\partial w_1}{\partial x} \\ &= \dots\end{aligned}$$

Dual Numbers

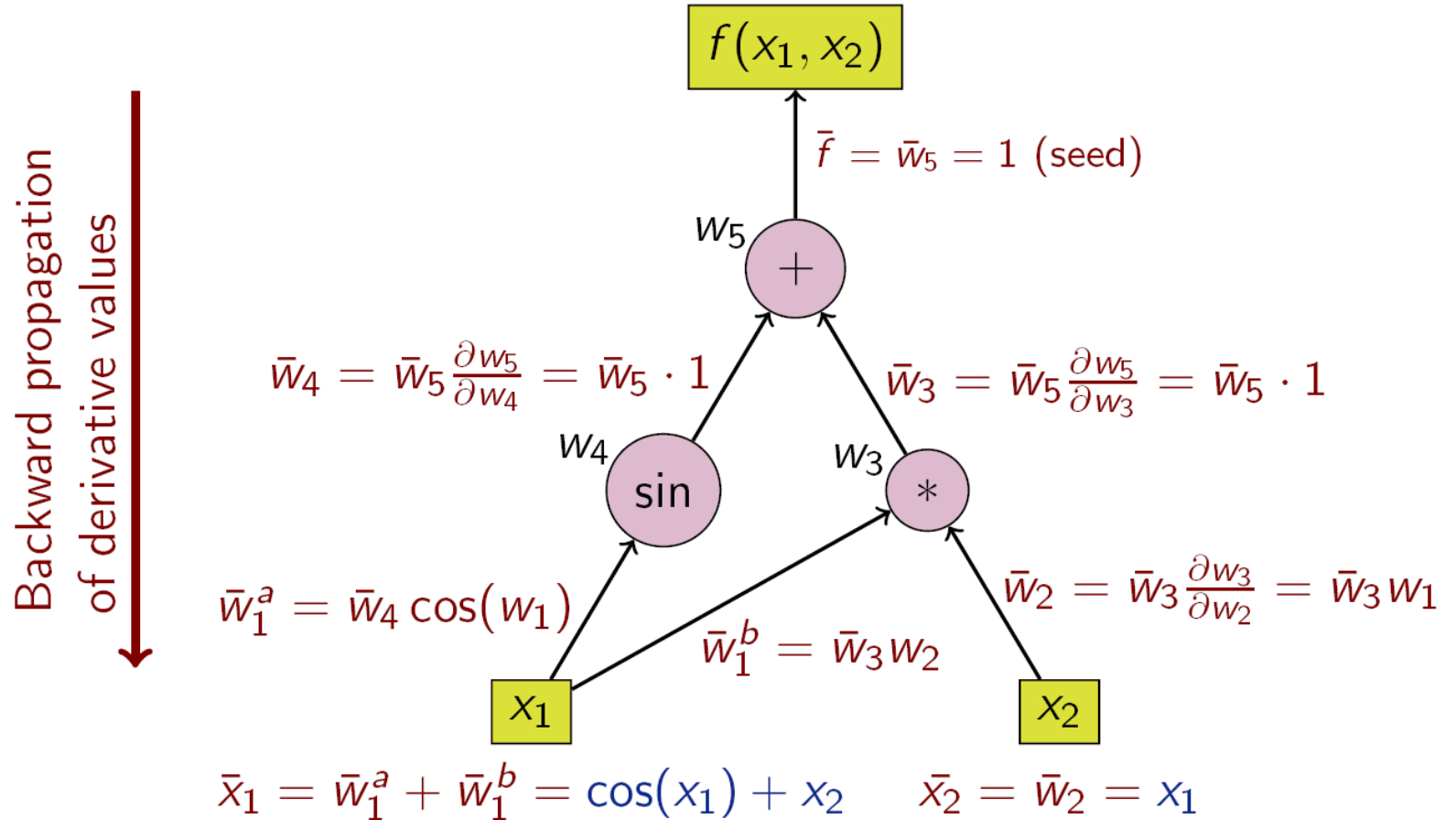
- Forward Mode
- **Derivatives**

Forward propagation
of derivative values



Dual Numbers

- Forward Mode
- **Derivatives**
- Reverse Mode
- **Adjoint**s



Naïve Reverse AD

- Source and Target Languages

Types:	$\sigma, \tau ::= \mathbb{R} \mid () \mid (\sigma, \tau) \mid \sigma \rightarrow \tau \mid \text{Int}$
Terms:	$s, t ::= x \mid () \mid (s, t) \mid \text{fst}(t) \mid \text{snd}(t) \mid s \ t \mid \lambda(x : \tau). t \mid \text{let } x : \tau = s \text{ in } t$ $\mid r$ (literal \mathbb{R} values) $\mid \text{op}(t_1, \dots, t_n)$ ($\text{op} \in \text{Op}_n$, primitive operation application ($\mathbb{R}^n \rightarrow \mathbb{R}$))

Fig. 4. The source language of all variants of this paper's reverse AD transformation. Int, the type of integers, is added as an example of a type that AD does not act upon.

Types:	$\sigma, \tau ::= \mathbb{R} \mid () \mid (\sigma, \tau) \mid \sigma \rightarrow \tau \mid \text{Int}$ $\mid \sigma \multimap \tau$ (linear functions)
Terms:	$s, t ::= x \mid () \mid (s, t) \mid \text{fst}(t) \mid \text{snd}(t) \mid s \ t \mid \lambda(x : \tau). t \mid \text{let } x : \tau = s \text{ in } t \mid r \mid \text{op}(t_1, \dots, t_n)$ $\mid \underline{\lambda}(z : \tau). b$ (linear lambda abstraction (τ a type without function arrows))
Linear function bodies:	
$b ::= () \mid (b, b') \mid \text{fst}(b) \mid \text{snd}(b)$	(tupling)
$\mid z$	(reference to $\underline{\lambda}$ -bound variable)
$\mid x \ b$	(linear function application; $x : \sigma \multimap \tau$ is an identifier)
$\mid \partial_i \text{op}(x_1, \dots, x_n)(b)$	($\text{op} \in \text{Op}_n$, i 'th partial derivative of op ($\mathbb{R}^n \rightarrow \mathbb{R}$))
$\mid b + b'$	(elementwise addition of results)
$\mid \underline{0}$	(zero of result type)

Fig. 5. The target language of the unoptimised variant of the reverse AD transformation. Components that are also in the source language (Fig. 4) are set in grey.

Naïve Reverse AD

- Source and Target Languages
- Linear Functions

$$f(0) = 0$$

$$f(x) + f(y) = f(x + y)$$

Types: $\sigma, \tau ::= \mathbb{R} \mid () \mid (\sigma, \tau) \mid \sigma \rightarrow \tau \mid \text{Int}$

Terms: $s, t ::= x \mid () \mid (s, t) \mid \text{fst}(t) \mid \text{snd}(t) \mid s \ t \mid \lambda(x : \tau). t \mid \text{let } x : \tau = s \text{ in } t$
 $\mid r$ (literal \mathbb{R} values)
 $\mid \text{op}(t_1, \dots, t_n)$ ($\text{op} \in \text{Op}_n$, primitive operation application ($\mathbb{R}^n \rightarrow \mathbb{R}$))

Fig. 4. The source language of all variants of this paper's reverse AD transformation. Int, the type of integers, is added as an example of a type that AD does not act upon.

Types: $\sigma, \tau ::= \mathbb{R} \mid () \mid (\sigma, \tau) \mid \sigma \rightarrow \tau \mid \text{Int}$
 $\mid \sigma \multimap \tau$ (linear functions)

Terms: $s, t ::= x \mid () \mid (s, t) \mid \text{fst}(t) \mid \text{snd}(t) \mid s \ t \mid \lambda(x : \tau). t \mid \text{let } x : \tau = s \text{ in } t \mid r \mid \text{op}(t_1, \dots, t_n)$
 $\mid \underline{\lambda}(z : \tau). b$ (linear lambda abstraction (τ a type without function arrows))

Linear function bodies:

$b ::= () \mid (b, b') \mid \text{fst}(b) \mid \text{snd}(b)$ (tupling)
 $\mid z$ (reference to $\underline{\lambda}$ -bound variable)
 $\mid x \ b$ (linear function application; $x : \sigma \multimap \tau$ is an identifier)
 $\mid \partial_i \text{op}(x_1, \dots, x_n)(b)$ ($\text{op} \in \text{Op}_n$, i 'th partial derivative of op ($\mathbb{R}^n \rightarrow \mathbb{R}$))
 $\mid b + b'$ (elementwise addition of results)
 $\mid \underline{0}$ (zero of result type)

Fig. 5. The target language of the unoptimised variant of the reverse AD transformation. Components that are also in the source language (Fig. 4) are set in grey.

Naïve Reverse AD

- Source and Target Languages
- Linear Functions
- Source Code Transformation

On types: $D_c^1[\mathbb{R}] = (\mathbb{R}, \mathbb{R} \multimap c)$ $D_c^1[()] = ()$ $D_c^1[(\sigma, \tau)] = (D_c^1[\sigma], D_c^1[\tau])$
 $D_c^1[\sigma \rightarrow \tau] = D_c^1[\sigma] \rightarrow D_c^1[\tau]$ $D_c^1[\text{Int}] = \text{Int}$
On environments: $D_c^1[\varepsilon] = \varepsilon$ $D_c^1[\Gamma, x : \tau] = D_c^1[\Gamma], x : D_c^1[\tau]$

On terms:
 If $\Gamma \vdash t : \tau$ then $D_c^1[\Gamma] \vdash D_c^1[t] : D_c^1[\tau]$
 $D_c^1[x : \tau] = x : D_c^1[\tau]$ $D_c^1[()] = ()$
 $D_c^1[(s, t)] = (D_c^1[s], D_c^1[t])$ $D_c^1[\text{fst}(t)] = \text{fst}(D_c^1[t])$
 $D_c^1[\text{snd}(t)] = \text{snd}(D_c^1[t])$ $D_c^1[s \ t] = D_c^1[s] \ D_c^1[t]$
 $D_c^1[\lambda(x : \tau). t] = \lambda(x : D_c^1[\tau]). D_c^1[t]$ $D_c^1[\text{let } x : \tau = s \text{ in } t] = \text{let } x : D_c^1[\tau] = D_c^1[s] \text{ in } D_c^1[t]$
 $D_c^1[r] = (r, \underline{\lambda}(z : \mathbb{R}). \underline{0})$
 $D_c^1[\text{op}(t_1, \dots, t_n)] = \text{let } (x_1, d_1) = D_c^1[t_1] \text{ in } \dots \text{ in let } (x_n, d_n) = D_c^1[t_n]$
 $\text{in } (\text{op}(x_1, \dots, x_n)$
 $\quad , \underline{\lambda}(z : \mathbb{R}). d_1 (\partial_1 \text{op}(x_1, \dots, x_n)(z)) + \dots + d_n (\partial_n \text{op}(x_1, \dots, x_n)(z)))$

Fig. 6. The naive code transformation from the source (Fig. 4) to the target (Fig. 5) language. The cases where D_c^1 just maps homomorphically over the source language are set in gray.

Naïve Reverse AD

- Source and Target Languages
- Linear Functions
- Source Code Transformation
- API Wrapper

$$\begin{aligned}
 \text{Interleave}_{\tau}^1 &: \forall c. (\tau, \tau \multimap c) \rightarrow \mathbf{D}_c^1[\tau] \\
 \text{Interleave}_{\mathbb{R}}^1 &= \lambda(x, d). (x, d) \\
 \text{Interleave}_{()}^1 &= \lambda(). ((), \underline{\lambda}(z : ().) . \underline{0}) \\
 \text{Interleave}_{(\sigma, \tau)}^1 &= \lambda((x, y), d). (\text{Interleave}_{\sigma}^1 (x, \underline{\lambda}(z : \sigma). d (z, \underline{0})), \text{Interleave}_{\tau}^1 (y, \underline{\lambda}(z : \tau). d (\underline{0}, z))) \\
 \text{Interleave}_{\text{Int}}^1 &= \lambda(n, d). n \\
 \text{Interleave}_{\sigma \rightarrow \tau}^1 &= \text{not defined!} \\
 \\
 \text{Deinterleave}_{\tau}^1 &: \forall c. \mathbf{D}_c^1[\tau] \rightarrow (\tau, \tau \multimap c) \\
 \text{Deinterleave}_{\mathbb{R}}^1 &= \lambda(x, d). (x, d) \\
 \text{Deinterleave}_{()}^1 &= \lambda(). ((), \underline{\lambda}(z : ().) . \underline{0}) \\
 \text{Deinterleave}_{(\sigma, \tau)}^1 &= \lambda(x, y). \text{let } (x_1, x_2) = \text{Deinterleave}_{\sigma}^1 x \\
 &\quad \text{in let } (y_1, y_2) = \text{Deinterleave}_{\tau}^1 y \\
 &\quad \text{in } ((x_1, y_1), \underline{\lambda}(z : (\sigma, \tau)). x_2 (\text{fst}(z)) + y_2 (\text{snd}(z))) \\
 \text{Deinterleave}_{\text{Int}}^1 &= \lambda n. (n, \underline{\lambda}(z : \text{Int}). \underline{0}) \\
 \text{Deinterleave}_{\sigma \rightarrow \tau}^1 &= \text{not defined!} \\
 \\
 \text{Wrap}^1 &: (\sigma \rightarrow \tau) \rightsquigarrow (\sigma \rightarrow (\tau, \tau \multimap \sigma)) \\
 \text{Wrap}^1[\lambda(x : \sigma). t] &= \lambda(x : \sigma). \text{let } x : \mathbf{D}_{\sigma}^1[\sigma] = \text{Interleave}_{\sigma}^1 (x, \text{id}) \text{ in Deinterleave}_{\tau}^1 (\mathbf{D}_{\sigma}^1[t])
 \end{aligned}$$

Fig. 7. Wrapper around \mathbf{D}_c^1 of Fig. 6.

Linear Factoring

- Duplicate backpropagator calls

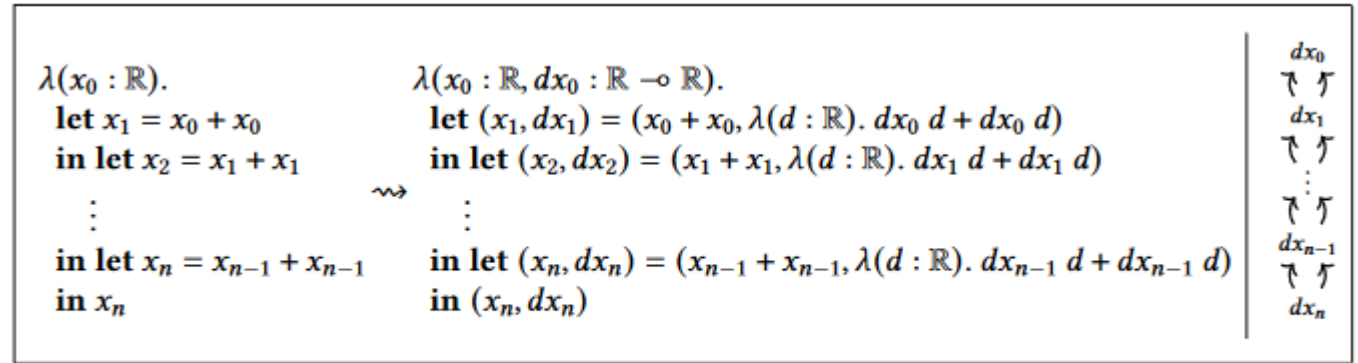


Fig. 2. Left: an example showing how naive dual-numbers reverse AD can result in exponential blow-up when applied to a program with sharing. Right: the dependency graph of the backpropagators dx_i .

Linear Factoring

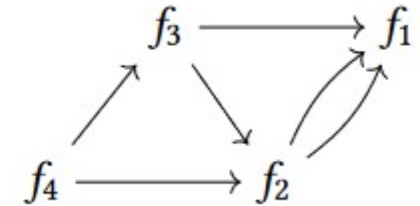
- Duplicate
backpropagator calls

$$f_1 = \lambda(z : \mathbb{R}). (0, (z, 0))$$

$$f_2 = \lambda(z : \mathbb{R}). f_1 (2 \cdot z) + f_1 (3 \cdot z)$$

$$f_3 = \lambda(z : \mathbb{R}). f_2 (4 \cdot z) + f_1 (5 \cdot z)$$

$$f_4 = \lambda(z : \mathbb{R}). f_2 z + f_3 (2 \cdot z)$$



Linear Factoring

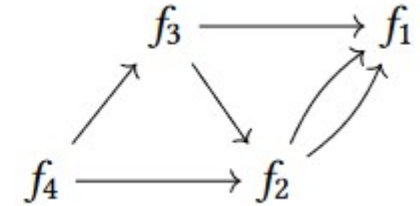
- Duplicate
backpropagator calls
- Linear Factoring

$$f_1 = \lambda(z : \mathbb{R}). (0, (z, 0))$$

$$f_2 = \lambda(z : \mathbb{R}). f_1 (2 \cdot z) + f_1 (3 \cdot z)$$

$$f_3 = \lambda(z : \mathbb{R}). f_2 (4 \cdot z) + f_1 (5 \cdot z)$$

$$f_4 = \lambda(z : \mathbb{R}). f_2 z + f_3 (2 \cdot z)$$



$$f(x) + f(y) = f(x + y)$$

Linear Factoring

- Duplicate
backpropagator calls
- Linear Factoring
- Execution Order

$\text{Staged } c = (c, \text{Map Int } (\mathbb{R} \multimap \text{Staged } c, \mathbb{R}))$

$\text{SCall} \quad : (\text{Int}, \mathbb{R} \multimap \text{Staged } c) \rightarrow \mathbb{R} \multimap \text{Staged } c$

$\text{SCall } (i, f) \ x = (\underline{0}, \{i \mapsto (f, x)\})$

Linear Factoring

- Duplicate
backpropagator calls
- Linear Factoring
- Execution Order
- Resolving “Staged c ”

```
SResolve ( $c : \sigma, m : \text{Map Int } (\mathbb{R} \multimap \text{Staged } \sigma, \mathbb{R})$ ) :=  
  if  $m$  is empty then  $c$   
  else let  $i = \text{highest key in } m$   
    in let  $(f, a) = \text{lookup } i \text{ in } m$   
    in let  $m' = \text{delete } i \text{ from } m$   
    in SResolve ( $f \ a +_{\text{Staged}} (c, m')$ )
```

Linear Factoring

- Duplicate
backpropagator calls
- Linear Factoring
- Execution Order
- Resolving “Staged c”
- New Transformation

On types:

$$\begin{aligned} D_c^2[\mathbb{R}] &= (\mathbb{R}, (\text{Int}, \mathbb{R} \multimap \text{Staged } c)) & D_c^2[()] &= () & D_c^2[(\sigma, \tau)] &= (D_c^2[\sigma], D_c^2[\tau]) \\ D_c^2[\sigma \rightarrow \tau] &= D_c^2[\sigma] \rightarrow \text{Int} \rightarrow (D_c^2[\tau], \text{Int}) & D_c^2[\text{Int}] &= \text{Int} \end{aligned}$$

On terms:

$$\text{If } \Gamma \vdash t : \tau \text{ then } D_c^2[\Gamma] \vdash D_c^2[t] : \text{Int} \rightarrow (D_c^2[\tau], \text{Int})$$

$$D_c^2[x : \tau] = \lambda i. (x : D_c^2[\tau], i)$$

$$D_c^2[(s, t)] = \lambda i. \text{let } (x, i') = D_c^2[s] \text{ } i \text{ in let } (y, i'') = D_c^2[t] \text{ } i' \text{ in } ((x, y), i'')$$

$$D_c^2[\text{let } x : \tau = s \text{ in } t] = \lambda i. \text{let } (x : D_c^2[\tau], i') = D_c^2[s] \text{ } i \text{ in } D_c^2[t] \text{ } i'$$

etc.

$$D_c^2[r] = \lambda i. ((r, (i, \underline{\lambda}(z : \mathbb{R}). 0_{\text{Staged}})), i + 1)$$

$$D_c^2[op(t_1, \dots, t_n)] =$$

$$\begin{aligned} &\lambda i. \text{let } ((x_1, d_1), i_1) = D_c^2[t_1] \text{ } i \text{ in } \dots \text{ in let } ((x_n, d_n), i_n) = D_c^2[t_n] \text{ } i_{n-1} \\ &\text{in } ((op(x_1, \dots, x_n), (i_n, \underline{\lambda}(z : \mathbb{R}). \text{SCall } d_1 (\partial_1 op(x_1, \dots, x_n)(z)) +_{\text{Staged}} \dots +_{\text{Staged}} \\ &\quad \text{SCall } d_n (\partial_n op(x_1, \dots, x_n)(z)))) \\ &\quad , i_n + 1) \end{aligned}$$

Changed wrapper:

$$\text{Wrap}^2 : (\sigma \rightarrow \tau) \rightsquigarrow (\sigma \rightarrow (\tau, \tau \multimap \sigma))$$

$$\begin{aligned} \text{Wrap}^2[\lambda(x : \sigma). t] &= \lambda(x : \sigma). \text{let } (x : D_\sigma^2[\sigma], i) = \text{Interleave}_\sigma^2(x, \text{SCotan}) 0 \\ &\quad \text{in let } (y, d) = \text{Deinterleave}_\tau^2(\text{fst}(D_\sigma^2[t] \text{ } i)) \\ &\quad \text{in } (y, \underline{\lambda}(z : \tau). \text{SResolve } (d \text{ } z)) \quad \text{— see main text for SResolve} \end{aligned}$$

$$\text{Interleave}_\tau^2 : \forall c. (\tau, \tau \multimap \text{Staged } c) \rightarrow \text{Int} \rightarrow (D_c^2[\tau], \text{Int})$$

$$\text{Interleave}_\mathbb{R}^2 = \lambda(x, d). \lambda i. ((x, (i, d)), i + 1)$$

$$\text{Interleave}_()^2 = \lambda((), d). \lambda i. ((), i)$$

$$\begin{aligned} \text{Interleave}_{(\sigma, \tau)}^2 &= \lambda((x, y), d). \lambda i. \text{let } (x', i') = \text{Interleave}_\sigma^2(x, \underline{\lambda}(z : \sigma). d \text{ } (z, 0)) \text{ } i \\ &\quad \text{in let } (y', i'') = \text{Interleave}_\tau^2(y, \underline{\lambda}(z : \tau). d \text{ } (0, z)) \text{ } i' \\ &\quad \text{in } ((x', y'), i'') \end{aligned}$$

$$\text{Interleave}_{\text{Int}}^2 = \lambda(n, d). \lambda i. (n, i)$$

Linear Factoring

- Duplicate
backpropagator calls
- Linear Factoring
- Execution Order
- Resolving “Staged c”
- New Transformation

etc.

$$D_c^2[r] = \lambda i. ((r, (i, \underline{\lambda}(z : \mathbb{R}). 0_{\text{Staged}})), i + 1)$$

$$D_c^2[op(t_1, \dots, t_n)] =$$

$$\lambda i. \text{let } ((x_1, d_1), i_1) = D_c^2[t_1] \text{ in } \dots \text{ in let } ((x_n, d_n), i_n) = D_c^2[t_n] \text{ in } i_{n-1}$$

$$\text{in } ((op(x_1, \dots, x_n), (i_n, \underline{\lambda}(z : \mathbb{R}). \text{SCall } d_1 (\partial_1 op(x_1, \dots, x_n)(z)) +_{\text{Staged}} \dots +_{\text{Staged}} \text{SCall } d_n (\partial_n op(x_1, \dots, x_n)(z))))$$

$$, i_n + 1)$$

Changed wrapper:

$$\text{Wrap}^2 : (\sigma \rightarrow \tau) \rightsquigarrow (\sigma \rightarrow (\tau, \tau \multimap \sigma))$$

$$\text{Wrap}^2[\lambda(x : \sigma). t] = \lambda(x : \sigma). \text{let } (x : D_\sigma^2[\sigma], i) = \text{Interleave}_\sigma^2(x, \text{SCotan}) \text{ in}$$

$$\text{let } (y, d) = \text{Deinterleave}_\tau^2(\text{fst}(D_\sigma^2[t] \ i)) \text{ in } (y, \underline{\lambda}(z : \tau). \text{SResolve } (d \ z)) \quad \text{— see main text for SResolve}$$

$$\text{Interleave}_\tau^2 : \forall c. (\tau, \tau \multimap \text{Staged } c) \rightarrow \text{Int} \rightarrow (D_c^2[\tau], \text{Int})$$

$$\text{Interleave}_\mathbb{R}^2 = \lambda(x, d). \lambda i. ((x, (i, d)), i + 1)$$

$$\text{Interleave}_()^2 = \lambda((), d). \lambda i. ((), i)$$

$$\text{Interleave}_{(\sigma, \tau)}^2 = \lambda((x, y), d). \lambda i. \text{let } (x', i') = \text{Interleave}_\sigma^2(x, \underline{\lambda}(z : \sigma). d \ (z, \underline{0})) \text{ in}$$

$$\text{let } (y', i'') = \text{Interleave}_\tau^2(y, \underline{\lambda}(z : \tau). d \ (\underline{0}, z)) \text{ in } ((x', y'), i'')$$

$$\text{Interleave}_{\text{Int}}^2 = \lambda(n, d). \lambda i. (n, i)$$

$\text{Deinterleave}_\tau^2$ gets type $\forall c. D_c^2[\tau] \rightarrow (\tau, \tau \multimap \text{Staged } c)$ and ignores the new Int in $D_c^2[\mathbb{R}]$.
 $\underline{0}$ changes to 0_{Staged} and $(+)$ changes to $(+_{\text{Staged}})$.

Staged interface:

$$\text{Staged } c = (c, \text{Map Int } (\mathbb{R} \multimap \text{Staged } c, \mathbb{R}))$$

$$0_{\text{Staged}} : \text{Staged } c \quad (+_{\text{Staged}}) : \text{Staged } c \rightarrow \text{Staged } c \rightarrow \text{Staged } c$$

$$0_{\text{Staged}} = (\underline{0}, \{\}) \quad (c, m) +_{\text{Staged}} (c', m') = (c + c', m \cup m') \quad \text{— with linear factoring}$$

$$\text{SCotan} : c \multimap \text{Staged } c \quad \text{SCall} : (\text{Int}, \mathbb{R} \multimap \text{Staged } c) \rightarrow \mathbb{R} \multimap \text{Staged } c$$

$$\text{SCotan } c = (c, \{\}) \quad \text{SCall } (i, f) \ x = (\underline{0}, \{i \mapsto (f, x)\})$$

Fig. 8. The monadically transformed code transformation (from Fig. 4 to Fig. 5 plus Staged operations), based on Fig. 6. Parts of D_c^2 and Interleave^2 that were simply lifted to monadic code are set in grey.

Current Complexity

- Primal Pass

$$\text{Wrap}^2[P] \ x$$

$$\text{Interleave}^2 : O(\text{size}(x))$$

$$\mathbf{D}_\sigma^2[P] : O(\text{cost}(P \ x))$$

$$\text{Deinterleave}^2 : O(\text{size}(P \ x))$$

$$\text{total} : O(\text{cost}(P \ x) + \text{size}(x))$$

Current Complexity

- Primal Pass
- Reverse Pass

$$\text{snd}(\text{Wrap}^2[P] \ x) \ dy$$

Current Complexity

- Primal Pass
- Reverse Pass
- **Expensive Monoid Operations**

$$O(\text{cost}(P\ x) + \text{size}(x))$$

$$\text{snd}(\text{Wrap}^2[P]\ x)\ dy$$

$$\begin{aligned} & (+_{\text{Staged}}) : \text{Staged } c \rightarrow \text{Staged } c \rightarrow \text{Staged } c \\ & (c, m) +_{\text{Staged}} (c', m') = (c + c', m \cup m') \quad \text{— with linear factoring} \end{aligned}$$

Current Complexity

- Primal Pass
- Reverse Pass
- **Expensive Monoid Operations**
- **Map operations in SResolve**

```
SResolve ( $c : \sigma, m : \text{Map Int } (\mathbb{R} \multimap \text{Staged } \sigma, \mathbb{R})$ ) :=  
  if  $m$  is empty then  $c$   
  else let  $i = \text{highest key in } m$   
    in let  $(f, a) = \text{lookup } i \text{ in } m$   
    in let  $m' = \text{delete } i \text{ from } m$   
    in SResolve ( $f \ a +_{\text{Staged}} (c, m')$ )
```

Current Complexity

- Primal Pass
- Reverse Pass
- **Expensive Monoid Operations**
- **Map operations in SResolve**
- **SCall in backpropagators**

$$(\underline{\lambda}(z : \mathbb{R}).\text{SCotan}(0, \dots, 0, z, 0, \dots, 0)) : O(\text{size}(x))$$

$$(\underline{\lambda}(z : \mathbb{R}).0_{\text{Staged}}) : O(\text{size}(x))$$

$$(\underline{\lambda}(z : \mathbb{R}).\text{SCall } d_1 (\partial_1 op(\dots)(z)) +_{\text{Staged}} \dots +_{\text{Staged}} \text{SCall } d_n (\partial_n op(\dots)(z))) : O(\text{size}(x))$$

Staged interface:

$$\text{Staged } c = (c, \text{Map Int } (\mathbb{R} \multimap \text{Staged } c, \mathbb{R}))$$

$$0_{\text{Staged}} : \text{Staged } c \quad (+_{\text{Staged}}) : \text{Staged } c \rightarrow \text{Staged } c \rightarrow \text{Staged } c$$

$$0_{\text{Staged}} = (\underline{0}, \{\}) \quad (c, m) +_{\text{Staged}} (c', m') = (c + c', m \cup m') \quad \text{— with linear factoring}$$

$$\text{SCotan} : c \multimap \text{Staged } c \quad \text{SCall} : (\text{Int}, \mathbb{R} \multimap \text{Staged } c) \rightarrow \mathbb{R} \multimap \text{Staged } c$$

$$\text{SCotan } c = (c, \{\}) \quad \text{SCall } (i, f) \ x = (\underline{0}, \{i \mapsto (f, x)\})$$

Cayley Transformation

- Monoid Isomorphism

$$(M, 0, +)$$

$$(M \rightarrow M, \text{id}, \circ)$$

Cayley Transformation

- Monoid Isomorphism
- Transforming Staged c

$$(M, 0, +)$$

$$(M \rightarrow M, \text{id}, \circ)$$

$$(\text{Staged } c, 0_{\text{Staged}}, +_{\text{Staged}})$$

$$(\text{Staged } c \rightarrow \text{Staged } c, \text{id}, \circ)$$

$$\text{Staged } c = (c, \text{Map Int } (\mathbb{R} \multimap (\text{Staged } c \rightarrow \text{Staged } c), \mathbb{R}))$$

$$\text{SCall} : (\text{Int}, \mathbb{R} \multimap (\text{Staged } c \rightarrow \text{Staged } c)) \rightarrow \mathbb{R} \multimap (\text{Staged } c \rightarrow \text{Staged } c)$$

$$\text{SCotan} : c \multimap (\text{Staged } c \rightarrow \text{Staged } c)$$

$$\text{SResolve} : (\text{Staged } c \rightarrow \text{Staged } c) \multimap c$$

$$\text{SCall } (i, f) \ x \ (c, m) = (c, \text{if } i \notin m \text{ then insert } i \mapsto (f, x) \text{ into } m \\ \text{else update } m \text{ at } i \text{ with } (\lambda(-, x'). (f, x + x'))))$$

$$\text{SCotan} : (c \rightarrow c) \multimap (\text{Staged } c \rightarrow \text{Staged } c)$$

$$\text{SCotan } f \ (c, m) = (f \ c, m)$$

Cayley Transformation

- Monoid Isomorphism
- Transforming Staged C
- New transformation

On types:

$$\begin{aligned} D_c^3[\mathbb{R}] &= (\mathbb{R}, (\text{Int}, \mathbb{R} \multimap (\text{Staged } c \rightarrow \text{Staged } c))) & D_c^3[()] &= () & D_c^3[(\sigma, \tau)] &= (D_c^3[\sigma], D_c^3[\tau]) \\ D_c^3[\sigma \rightarrow \tau] &= D_c^3[\sigma] \rightarrow \text{Int} \rightarrow (D_c^3[\tau], \text{Int}) & D_c^3[\text{Int}] &= \text{Int} \end{aligned}$$

On terms:

$$\text{If } \Gamma \vdash t : \tau \text{ then } D_c^3[\Gamma] \vdash D_c^3[t] : \text{Int} \rightarrow (D_c^3[\tau], \text{Int})$$

Same as D_c^2 , except with 'id' in place of 0_{Staged} and 'o' in place of $(+_{\text{Staged}})$.

Changed wrapper:

$$\text{Wrap}^3 : (\sigma \rightarrow \tau) \rightsquigarrow (\sigma \rightarrow (\tau, \tau \multimap \sigma))$$

$$\begin{aligned} \text{Wrap}^3[\lambda(x : \sigma). t] &= \lambda(x : D_\sigma^3[\sigma], i) = \text{Interleave}_\sigma^3(x, \text{SCotan}) 0 \\ &\quad \text{in let } (y, d) = \text{Deinterleave}_\tau^3(\text{fst}(D_\sigma^3[t] i)) \\ &\quad \text{in } (y, \underline{\lambda}(z : \tau). \text{SResolve } (d z)) \end{aligned}$$

$$\text{Interleave}_\tau^3 : \forall c. (\tau, (\tau \rightarrow \tau) \multimap (\text{Staged } c \rightarrow \text{Staged } c)) \rightarrow \text{Int} \rightarrow (D_c^3[\tau], \text{Int})$$

$$\begin{aligned} \text{Interleave}_\mathbb{R}^3 &= \lambda(x, d). \lambda i. ((x, (i, \underline{\lambda}(z : \mathbb{R}). d (\lambda(a : \mathbb{R}). z + a))), \\ &\quad , i + 1) \end{aligned}$$

$$\text{Interleave}_{()}^3 = \lambda((), d). \lambda i. ((), i)$$

$$\text{Interleave}_{(\sigma, \tau)}^3 = \lambda((x, y), d). \lambda i.$$

$$\begin{aligned} &\quad \text{let } (x', i') = \text{Interleave}_\sigma^3(x, \lambda(f : \sigma \rightarrow \sigma). d (\lambda((v, w) : (\sigma, \tau)). (f v, w))) i \\ &\quad \text{in let } (y', i'') = \text{Interleave}_\tau^3(y, \lambda(f : \tau \rightarrow \tau). d (\lambda((v, w) : (\sigma, \tau)). (v, f w))) i' \\ &\quad \text{in } ((x', y'), i'') \end{aligned}$$

$$\text{Interleave}_{\text{Int}}^3 = \lambda(n, d). \lambda i. (n, i)$$

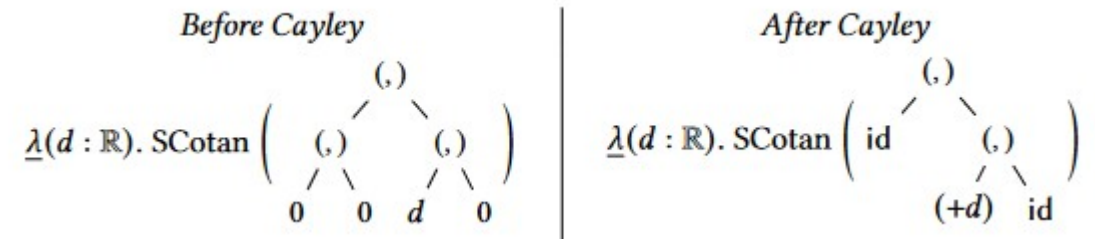
$$\text{Deinterleave}_\tau^3 : \forall c. D_c^3[\tau] \rightarrow (\tau, \tau \multimap (\text{Staged } c \rightarrow \text{Staged } c))$$

(Same as Deinterleave^2 in Fig. 8, except with id and (o) in place of 0_{Staged} and $(+_{\text{Staged}})$)

Fig. 9. The Cayley-transformed code transformation, based on Fig. 8. Grey parts are unchanged.

Cayley Transformation

- Monoid Isomorphism
- Transforming Staged C
- New transformation
- Remaining Problem



Efficient Gradient Updates

- Only collect scalars

$$\text{Staged } c = (c, \text{Map Int } (\mathbb{R} \multimap (\text{Staged } c \rightarrow \text{Staged } c), \mathbb{R}))$$

$$\text{Staged } c = (\text{Map Int } \mathbb{R}, \text{Map Int } (\mathbb{R} \multimap (\text{Staged } c \rightarrow \text{Staged } c), \mathbb{R}))$$

Efficient Gradient Updates

- Only collect scalars
- Remaining problems:
- **Expensive Monoid Operations**
- **Map Operations in SResolve**
- **SCall in backpropagators**

$$O((\text{cost}(P\ x) + \text{size}(x)) \log(\text{cost}(P\ x) + \text{size}(x)))$$

Efficient Gradient Updates

- Only collect scalars
- Remaining problems:
- ~~Expensive Monoid~~
Operations
- **Map Operations in SResolve**
- ~~SCall in backpropagators~~

$$O((\text{cost}(P\ x) + \text{size}(x)) \log(\text{cost}(P\ x) + \text{size}(x)))$$

Mutable Arrays

- Map operations in SResolve

```
SResolve ( $c : \sigma, m : \text{Map Int } (\mathbb{R} \multimap \text{Staged } \sigma, \mathbb{R})$ ) :=  
  if  $m$  is empty then  $c$   
  else let  $i = \text{highest key in } m$   
    in let  $(f, a) = \text{lookup } i \text{ in } m$   
    in let  $m' = \text{delete } i \text{ from } m$   
    in SResolve ( $f \ a +_{\text{Staged}} (c, m')$ )
```

Mutable Arrays

- Map operations in SResolve
- Maps to Arrays

$$\text{Staged } c = (\text{Map Int } \mathbb{R}, \text{Map Int } (\mathbb{R} \multimap (\text{Staged } c \rightarrow \text{Staged } c), \mathbb{R}))$$
$$\text{Staged } c = (\text{Array } \mathbb{R}, \text{Array } (\mathbb{R} \multimap (\text{Staged } c \xrightarrow{R} \text{Staged } c), \mathbb{R}))$$

More Optimizations

- No cotangent collection array

$$(\text{Array } \mathbb{R}, \text{Array } (\mathbb{R} \multimap (\text{Staged } c \xrightarrow{R} \text{Staged } c), \mathbb{R}))$$
$$\text{Array } (\mathbb{R} \multimap (\text{Staged } c \xrightarrow{R} \text{Staged } c), \mathbb{R})$$

More Optimizations

- No cotangent collection array
- Defunctionalization of backpropagators

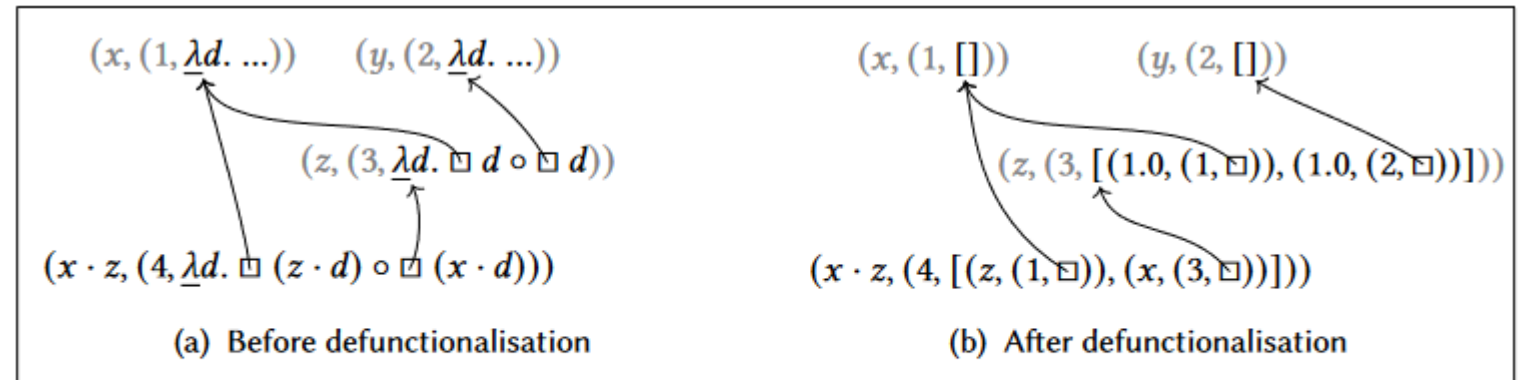


Fig. 10. The sharing structure before and after defunctionalisation. SCall is elided here; in Fig. 10a, the backpropagator calls are depicted as if they are still normal calls. Boxes (\square) are the same in-memory value as the value their arrow points to; two boxes pointing to the same value indicates that this value is *shared*: referenced in two places.

More Optimizations

- No cotangent collection array
- Defunctionalization of backpropagators
- Taping

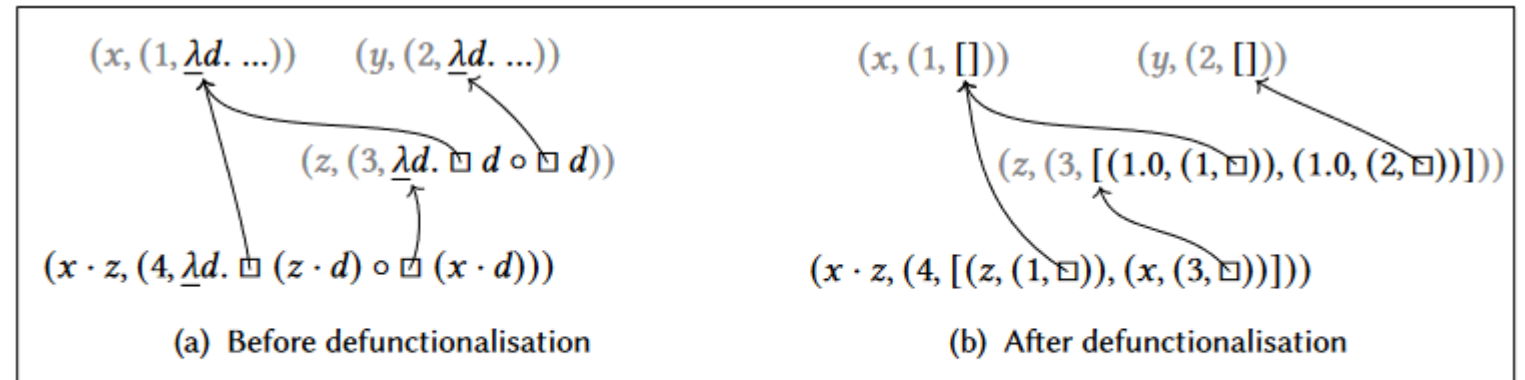


Fig. 10. The sharing structure before and after defunctionalisation. SCall is elided here; in Fig. 10a, the backpropagator calls are depicted as if they are still normal calls. Boxes (□) are the same in-memory value as the value their arrow points to; two boxes pointing to the same value indicates that this value is *shared*: referenced in two places.

Summary

- Dual-Numbers Reverse AD
- Linear Factoring
- Cayley Transformation
- Efficient Gradient Updates
- Mutable Arrays
- Even more optimizations

Summary

- Dual-Numbers Reverse AD
- Linear Factoring
- Cayley Transformation
- Efficient Gradient Updates
- Mutable Arrays
- Even more optimizations
- Extensions to the source language
- Implementation and Performance

Questions?