

# AD for an Array Language Supporting Nested Parallelism

Cosmin E. Oancea ([cosmin.oancea@diku.dk](mailto:cosmin.oancea@diku.dk))

in collaboration with

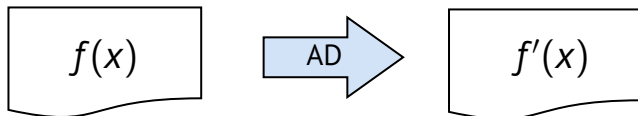
Troels Henriksen, Ola Rønning and Robert Schenck

Department of Computer Science (DIKU)  
University of Copenhagen

7th of June 2023, MIAPbP Workshop

# Overview

- **Automatic differentiation (AD)** is a program transformation for differentiation.



- This talk presents an AD technique [1] for a functional, **high-level**, and **nested-parallel** array language (Futhark).
- All parallelism is made explicit via **parallel combinators**—map, reduce, scan (prefix sum), scatter, etc.

[1] Robert Schenck, Ola Rønning, Troels Henriksen and Cosmin Oancea, "AD for an Array Language with Nested Parallelism", In Procs of SC22: International Conference for High Performance Computing, Networking, Storage and Analysis, 2022.

Motivating Example for AD

Gentle Introduction to AD

Key Idea #1: Redundant Execution Instead of Tape

Key Idea #2: Parallel Constructs are Differentiated at a High Level

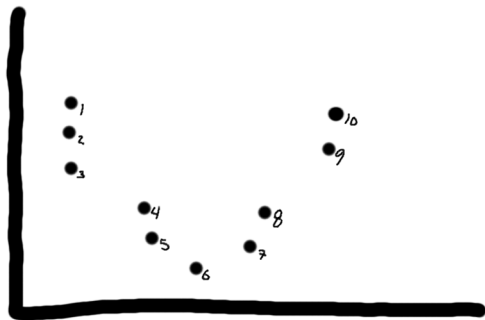
Key Idea #3: Translate Accumulators to Specialized Constructs

Experimental Results: Competitive with State of the Art

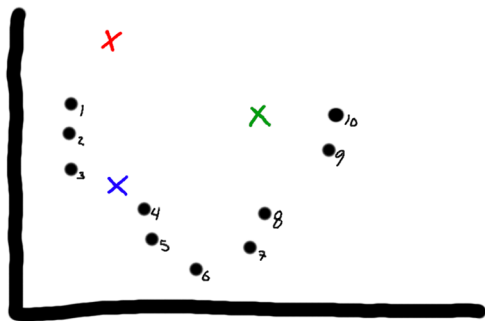
Demo

Extra Slides

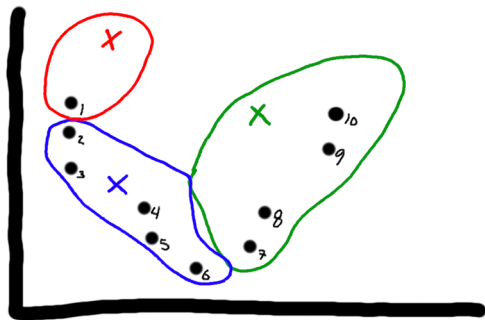
# Illustrative Example: $k$ -means clustering in 2D



# Illustrative Example: $k$ -means clustering in 2D



# Illustrative Example: $k$ -means clustering in 2D



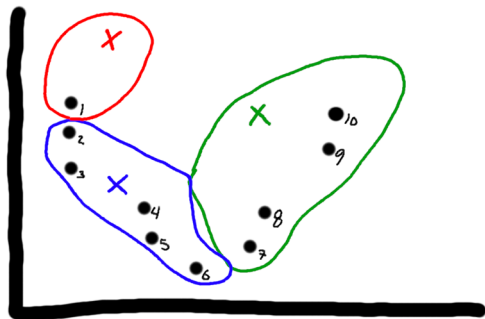
# Illustrative Example: $k$ -means clustering in 2D

cluster sizes

$$|1| = 1$$

$$|1+1+1+1| = 4$$

$$|1+1+1+1+1| = 5$$



# Illustrative Example: $k$ -means clustering in 2D

Cluster sizes

$$|1| = 1$$

$$|1+1+1+1| = 4$$

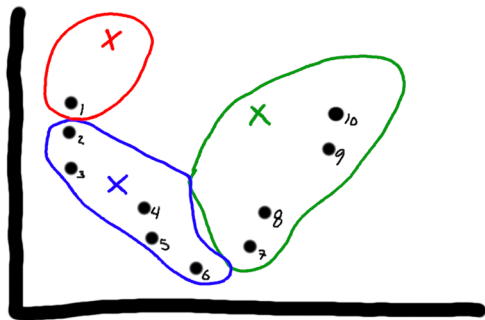
$$|1+1+1+1+1| = 5$$

Cluster sums

$$| \bullet_1 | = +$$

$$| \bullet_7 + \bullet_8 + \bullet_9 + \bullet_{10} | = +$$

$$| \bullet_2 + \bullet_3 + \bullet_4 + \bullet_5 + \bullet_6 | = +$$





# Illustrative Example: $k$ -means clustering in 2D

Cluster sizes

$$|1| = 1$$

$$|1+1+1+1| = 4$$

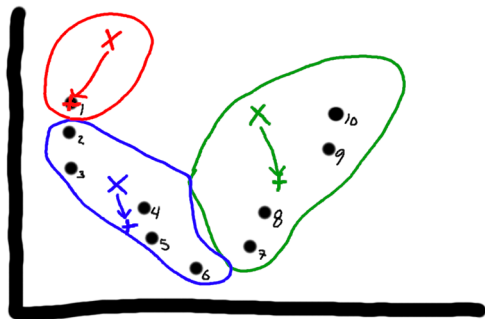
$$|1+1+1+1+1| = 5$$

Cluster sums

$$| \bullet_1 | = +$$

$$| \bullet_7 + \bullet_8 + \bullet_9 + \bullet_{10} | = +$$

$$| \bullet_2 + \bullet_3 + \bullet_4 + \bullet_5 + \bullet_6 | = +$$



# Illustrative Example: $k$ -means clustering in 2D

Cluster sizes

$$| 1 = 1$$

$$| 1+1+1+1 = 4$$

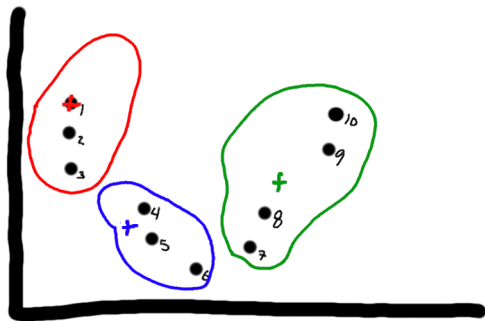
$$| 1+1+1+1+1 = 5$$

Cluster sums

$$| \bullet_1 = +$$

$$| \bullet_7 + \bullet_8 + \bullet_9 + \bullet_{10} = +$$

$$| \bullet_2 + \bullet_3 + \bullet_4 + \bullet_5 + \bullet_6 = +$$



# Illustrative Example: $k$ -means clustering in 2D

## Cluster sizes

$$|1| = 1$$

$$|1+1+1+1| = 4$$

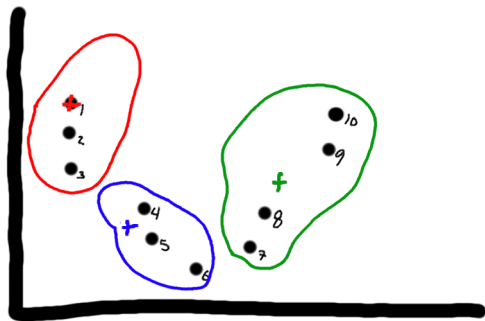
$$|1+1+1+1+1| = 5$$

## Cluster sums

$$| \bullet_1 | = +$$

$$| \bullet_7 + \bullet_8 + \bullet_9 + \bullet_{10} | = +$$

$$| \bullet_2 + \bullet_3 + \bullet_4 + \bullet_5 + \bullet_6 | = +$$



Generalized histograms allow a two-slide efficient implementation.

[2] Troels Henriksen, Sune Hellfritzsch, P. Sadayappan and Cosmin Oancea, "Compiling Generalized Histograms for GPU", In Procs of SC20: International Conference for High Performance Computing, Networking, Storage and Analysis, 2020.

# Mathematical Formulation of $k$ -means clustering

Given a set of  $n$  points  $P$  in a  $d$ -dimensional space, one must find the  $k$  points  $C$  that minimize the cost function:

$$f(C) = \sum_{p \in P} \min \{ \|p - c\|^2, c \in C \}$$

# Mathematical Formulation of $k$ -means clustering

Given a set of  $n$  points  $P$  in a  $d$ -dimensional space, one must find the  $k$  points  $C$  that minimize the cost function:

$$f(C) = \sum_{p \in P} \min \{ \|p - c\|^2, c \in C \}$$

This problem can be solved by applying Newton's Method, i.e.,

$$C_{i+1} = C_i - \nabla f(C_i) \cdot H_f(C_i)^{-1}$$

# Motivation for AD

- Initially: because we heard that AD is a challenging code transformation, and we are always up to a good challenge.
- Because AD enables simple and generic optimization recipes
  - ▶ Gradient descent and Newton's Method
  - ▶ In Probabilistic ML:
    - ▶ Hamiltonian MC,
    - ▶ Newtonian MC,
    - ▶ variational inference,
    - ▶ ...

**AD: worthy to be supported as a first-class citizen of parallel languages?**

Motivating Example for AD

## Gentle Introduction to AD

Key Idea #1: Redundant Execution Instead of Tape

Key Idea #2: Parallel Constructs are Differentiated at a High Level

Key Idea #3: Translate Accumulators to Specialized Constructs

Experimental Results: Competitive with State of the Art

Demo

Extra Slides

# Math: Simple Differentiation Rules $h : \mathbb{R} \rightarrow \mathbb{R}$

Basic Rules:

sin:  $\frac{\partial \sin x}{\partial x} = \cos x$

cos:  $\frac{\partial \cos x}{\partial x} = -\sin x$

pow:  $\frac{\partial x^r}{\partial x} = r \cdot x^{r-1}, \forall r \in \mathbb{R}, r \neq 0$

log:  $\frac{\partial \log_c x}{\partial x} = \frac{1}{x \cdot \ln c}$  where

$\ln x = \log_e x, e = 2.718281828459\dots$  Euler's number.

Linear: If  $h(x) = f(x) + g(x)$  then  $\frac{\partial h}{\partial x} = \frac{\partial f}{\partial x} + \frac{\partial g}{\partial x}$

...

Product: If  $h(x) = f(x) \cdot g(x)$  then  $\frac{\partial h}{\partial x} = \frac{\partial f}{\partial x} \cdot g + f \cdot \frac{\partial g}{\partial x}$

Power: If  $h(x) = f(x)^r$  then  $\frac{\partial h}{\partial x} = n \cdot f^{r-1} \cdot \frac{\partial f}{\partial x}$

...



# Math: Simple Differentiation Rules $h : \mathbb{R} \rightarrow \mathbb{R}$

Basic Rules:

$$\text{sin: } \frac{\partial \sin x}{\partial x} = \cos x$$

$$\text{cos: } \frac{\partial \cos x}{\partial x} = -\sin x$$

$$\text{pow: } \frac{\partial x^r}{\partial x} = r \cdot x^{r-1}, \quad \forall r \in \mathbb{R}, r \neq 0$$

$$\text{log: } \frac{\partial \log_c x}{\partial x} = \frac{1}{x \cdot \ln c} \text{ where}$$

$\ln x = \log_e x$ ,  $e = 2.718281828459\dots$  Euler's number.

$$\text{Linear: If } h(x) = f(x) + g(x) \text{ then } \frac{\partial h}{\partial x} = \frac{\partial f}{\partial x} + \frac{\partial g}{\partial x}$$

...

$$\text{Product: If } h(x) = f(x) \cdot g(x) \text{ then } \frac{\partial h}{\partial x} = \frac{\partial f}{\partial x} \cdot g + f \cdot \frac{\partial g}{\partial x}$$

$$\text{Power: If } h(x) = f(x)^r \text{ then } \frac{\partial h}{\partial x} = n \cdot f^{r-1} \cdot \frac{\partial f}{\partial x}$$

...

**The Chain Rule (the rule to bind them all):**

If  $h(x) = (f \circ g)(x) = f(g(x))$  then

$$\frac{\partial h(x)}{\partial x} = \frac{\partial f(z)|_{z=g(x)}}{\partial z} \cdot \frac{\partial g(x)}{\partial x}$$

# Math: Chain Rule for Multi-Variate Functions

Given  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ ,  $f(x) = (f_1(x), \dots, f_m(x))$ ,  $f_i : \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $a \in \mathbb{R}^n$

$$J_f(a) = \left[ \frac{\partial f(a)}{\partial a_1}(a_1) \dots \frac{\partial f(a)}{\partial a_n}(a_n) \right] = \begin{bmatrix} \frac{\partial f_1(a)}{\partial a_1}(a_1) & \dots & \frac{\partial f_1(a)}{\partial a_n}(a_n) \\ \frac{\partial f_m(a)}{\partial a_1}(a_1) & \dots & \frac{\partial f_m(a)}{\partial a_n}(a_n) \end{bmatrix}$$

$$\frac{\partial f(a)}{\partial a}(y) = J_f(a) \cdot y$$

Chain Rule for  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ ,  $g : \mathbb{R}^m \rightarrow \mathbb{R}^k$  and  $p \in \mathbb{R}^n$  is:

$$J_{g \circ f}(p) = J_g(f(p)) \cdot J_f(p)$$

# Math Intuition for Differentiating a Program

Consider  $f_i : \mathbb{R}^{n_{i-1}} \rightarrow \mathbb{R}^{n_i}, \forall i = 1 \dots m$ , and “program”  $P$  defined as:

$$P : \mathbb{R}^{n_0} \rightarrow \mathbb{R}^{n_m} = f_m \circ f_{m-1} \circ \dots \circ f_2 \circ f_1$$

Applying the Chain Rule for some  $x \in \mathbb{R}^{n_0}$  results in:

$$J_P(x) = J_{f_m}(f_{m-1}(\dots f_1(x))) \cdot \dots \cdot J_{f_2}(f_1(x)) \cdot J_{f_1}(x)$$

**Forward vs Reverse-Mode AD answers the question:**

**Should Jacobians be multiplied from right-to-left or left-to-right?**

# Intuition: Reverse-Mode AD

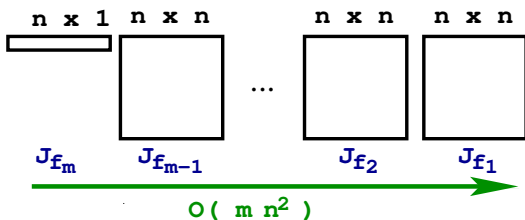
Consider  $f_i : \mathbb{R}^{n_{i-1}} \rightarrow \mathbb{R}^{n_i}, \forall i = 1 \dots m$ , and “program”  $P$  defined as:

$$P : \mathbb{R}^{n_0} \rightarrow \mathbb{R}^{n_m} = f_m \circ f_{m-1} \circ \dots \circ f_2 \circ f_1$$

Applying the Chain Rule for some  $x \in \mathbb{R}^{n_0}$  results in:

$$J_P(x) = J_{f_m}(f_{m-1}(\dots f_1(x))) \cdot \dots \cdot J_{f_2}(f_1(x)) \cdot J_{f_1}(x)$$

**Reverse-Mode AD multiplies Jacobians from left to right.**



# Intuition: Reverse-Mode AD

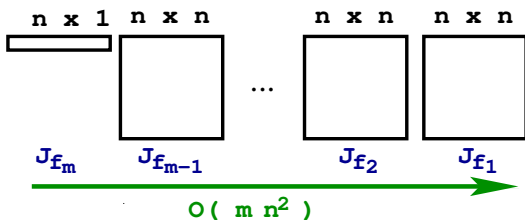
Consider  $f_i : \mathbb{R}^{n_{i-1}} \rightarrow \mathbb{R}^{n_i}, \forall i = 1 \dots m$ , and “program”  $P$  defined as:

$$P : \mathbb{R}^{n_0} \rightarrow \mathbb{R}^{n_m} = f_m \circ f_{m-1} \circ \dots \circ f_2 \circ f_1$$

Applying the Chain Rule for some  $x \in \mathbb{R}^{n_0}$  results in:

$$J_P(x) = J_{f_m}(f_{m-1}(\dots f_1(x))) \cdot \dots \cdot J_{f_2}(f_1(x)) \cdot J_{f_1}(x)$$

**Reverse-Mode AD multiplies Jacobians from left to right.**



- Preferred when the input size is larger than the size of the result (ML case);
- The *primal trace* (original program) is first executed to save all intermediate program values on a *tape*. The tape is subsequently used by the *return sweep*, which computes the *adjoint* of each variable in reverse program order.

## Reverse AD: Core Re-Write Rule for Scalar Code

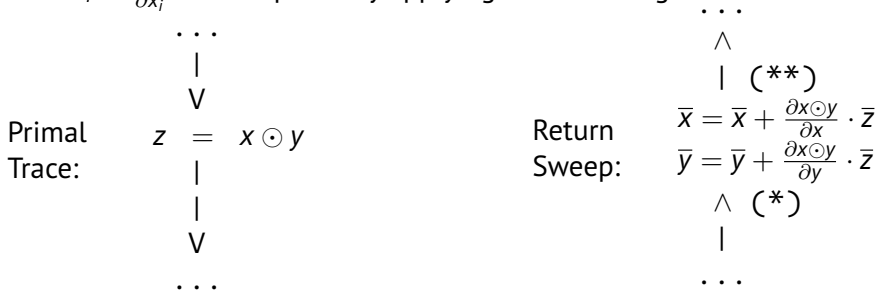
For program  $P(\dots x_i \dots) = y \in \mathbb{R}$ , reverse AD computes the **adjoint** of each (intermediate) program variable  $t$ , denoted  $\bar{t} = \frac{\partial y}{\partial t}$ , i.e., the sensitivity of the output to changes in  $t$ .

# Reverse AD: Core Re-Write Rule for Scalar Code

For program  $P(\dots x_i \dots) = y \in \mathbb{R}$ , reverse AD computes the **adjoint** of each (intermediate) program variable  $t$ , denoted  $\bar{t} = \frac{\partial y}{\partial t}$ , i.e., the sensitivity of the output to changes in  $t$ .

Initially  $\bar{y} = \frac{\partial y}{\partial y} = 1$ , and eventually the adjoints of the input

$\bar{x}_i = \frac{\partial y}{\partial x_i}$  are computed by applying the following re-write rule:



- (\*) final value of  $\bar{z}$  is known; the values of  $x$  and  $y$  need to be available, e.g., used in  $\frac{\partial x \odot y}{\partial x}$
- (\*\*)  $\bar{z}$  is dead after this point;  $\bar{x}$  and  $\bar{y}$  are still under computation;
  - adjoints are initialized (zeroed) before their first use.

# Illustrating Reverse AD on straightline scalar code

**Adjoint** of var  $t$  is  $\bar{t} \equiv \frac{\partial y}{\partial t}$ . **Goal:** compute the adjoints of the input.

**Core Rewrite Rule:**  $z = f(x, y) \Rightarrow$

$$\bar{x}_+ = \frac{\partial f(x,y)}{\partial x} \cdot \bar{z}$$
$$\bar{y}_+ = \frac{\partial f(x,y)}{\partial y} \cdot \bar{z}$$

$P(x_0, x_1) :$   
 $t_0 = \sin(x_0)$   
 $t_1 = x_1 \cdot t_0$   
 $y = x_0 + t_1$   
**return**  $y$

$\Rightarrow$

$P'(x_0, x_1) :$   
 $t_0 = \sin(x_0)$   
 $t_1 = x_1 \cdot t_0$   
 $y = x_0 + t_1$



# Illustrating Reverse AD on straightline scalar code

**Adjoint** of var  $t$  is  $\bar{t} \equiv \frac{\partial y}{\partial t}$ . **Goal:** compute the adjoints of the input.

**Core Rewrite Rule:**  $y = f(x_1, x_2) \Rightarrow$

$$\bar{x}_1 + = \frac{\partial f(x_1, x_2)}{\partial x_1} \cdot \bar{z}$$
$$\bar{x}_2 + = \frac{\partial f(x_1, x_2)}{\partial x_2} \cdot \bar{z}$$

$P(x_0, x_1) :$   
 $t_0 = \sin(x_0)$   
 $t_1 = x_1 \cdot t_0$   
 $y = x_0 + t_1$   
**return**  $y$

$\Rightarrow$

$P'(x_0, x_1) :$   
 $t_0 = \sin(x_0)$   
 $t_1 = x_1 \cdot t_0$   
 $y = x_0 + t_1$   
 $\bar{y} = 1; \quad \bar{x}_0 = 0; \quad \bar{t}_1 = 0$   
 $\bar{x}_0 = \bar{x}_0 + \frac{\partial(x_0+t_1)}{\partial x_0} \cdot \bar{y} = 1 \cdot 1 = 1$   
 $\bar{t}_0 = \bar{t}_0 + \frac{\partial(x_0+t_1)}{\partial t_1} \cdot \bar{y} = 1 \cdot 1 = 1$

# Illustrating Reverse AD on straightline scalar code

**Adjoint** of var  $t$  is  $\bar{t} \equiv \frac{\partial y}{\partial t}$ . **Goal:** compute the adjoints of the input.

**Core Rewrite Rule:**  $y = f(x_1, x_2) \Rightarrow$

$$\bar{x}_1 + = \frac{\partial f(x_1, x_2)}{\partial x_1} \cdot \bar{z}$$
$$\bar{x}_2 + = \frac{\partial f(x_1, x_2)}{\partial x_2} \cdot \bar{z}$$

$P(x_0, x_1) :$		$P'(x_0, x_1) :$
$t_0 = \sin(x_0)$	$\Rightarrow$	$t_0 = \sin(x_0)$
$t_1 = x_1 \cdot t_0$		$t_1 = x_1 \cdot t_0$
$y = x_0 + t_1$		$y = x_0 + t_1$
<b>return</b> $y$		$\bar{y} = 1; \quad \bar{x}_0 = 0; \quad \bar{t}_1 = 0$
		$\bar{x}_0 = \bar{x}_0 + \frac{\partial(x_0+t_1)}{\partial x_0} \cdot \bar{y} = 1 \cdot 1 = 1$
		$\bar{t}_1 = \bar{t}_0 + \frac{\partial(x_0+t_1)}{\partial t_1} \cdot \bar{y} = 1 \cdot 1 = 1$
		$\bar{x}_1 = 0; \quad \bar{t}_0 = 0$
		$\bar{x}_1 = \bar{x}_1 + \frac{\partial(x_1 \cdot t_0)}{\partial x_1} \cdot \bar{t}_1 = \bar{x}_1 + t_0 \cdot \bar{t}_1$
		$\bar{t}_0 = \bar{t}_0 + \frac{\partial(x_1 \cdot t_0)}{\partial t_0} \cdot \bar{t}_1 = \bar{t}_0 + x_1 \cdot \bar{t}_1$

# Illustrating Reverse AD on straightline scalar code

**Adjoint** of var  $t$  is  $\bar{t} \equiv \frac{\partial y}{\partial t}$ . **Goal:** compute the adjoints of the input.

**Core Rewrite Rule:**  $y = f(x_1, x_2) \Rightarrow$

$$\bar{x}_1 + = \frac{\partial f(x_1, x_2)}{\partial x_1} \cdot \bar{z}$$
$$\bar{x}_2 + = \frac{\partial f(x_1, x_2)}{\partial x_2} \cdot \bar{z}$$

```
P(x0, x1) :  
  t0 = sin(x0)  
  t1 = x1 · t0  
  y = x0 + t1  
  return y
```

$\Rightarrow$

```
P'(x0, x1) :  
  t0 = sin(x0)  
  t1 = x1 · t0  
  y = x0 + t1  
   $\bar{y} = 1; \quad \bar{x}_0 = 0; \quad \bar{t}_1 = 0$   
   $\bar{x}_0 = \bar{x}_0 + \frac{\partial(x_0 + t_1)}{\partial x_0} \cdot \bar{y} = 1 \cdot 1 = 1$   
   $\bar{t}_1 = \bar{t}_0 + \frac{\partial(x_0 + t_1)}{\partial t_1} \cdot \bar{y} = 1 \cdot 1 = 1$   
   $\bar{x}_1 = 0; \quad \bar{t}_0 = 0$   
   $\bar{x}_1 = \bar{x}_1 + \frac{\partial(x_1 \cdot t_0)}{\partial x_1} \cdot \bar{t}_1 = \bar{x}_1 + t_0 \cdot \bar{t}_1$   
   $\bar{t}_0 = \bar{t}_0 + \frac{\partial(x_1 \cdot t_0)}{\partial t_0} \cdot \bar{t}_1 = \bar{t}_0 + x_1 \cdot \bar{t}_1$   
   $\bar{x}_0 = \bar{x}_0 + \frac{\partial \sin(x_0)}{\partial x_0} \cdot \bar{t}_0 = \bar{x}_0 + \cos(x_0) \cdot \bar{t}_0$   
  return  $\bar{x}_0, \bar{x}_1$ 
```

# Classical API for AD

## Reverse-mode:

$$\text{vjp} : (f : \alpha \rightarrow \beta) \rightarrow (x : \alpha) \rightarrow (\bar{y} : \beta) \rightarrow \alpha$$

## Forward mode:

$$\text{jvp} : (f : \alpha \rightarrow \beta) \rightarrow (x : \alpha) \rightarrow (dx : \beta) \rightarrow \beta$$

# Classical API for AD

## Reverse-mode:

$$\text{vjp} : (f : \alpha \rightarrow \beta) \rightarrow (x : \alpha) \rightarrow (\bar{y} : \beta) \rightarrow \alpha$$

## Forward mode:

$$\text{jvp} : (f : \alpha \rightarrow \beta) \rightarrow (x : \alpha) \rightarrow (dx : \beta) \rightarrow \beta$$

Matrix Multiplication in Futhark (C = A\*B):

```
let C =  
  map (\A_row ->  
    map (\B_col ->  
      map2 (*) A_row B_col  
      |> reduce (+) 0  
    ) (transpose B)  
  ) A
```

Motivating Example for AD

Gentle Introduction to AD

**Key Idea #1: Redundant Execution Instead of Tape**

Key Idea #2: Parallel Constructs are Differentiated at a High Level

Key Idea #3: Translate Accumulators to Specialized Constructs

Experimental Results: Competitive with State of the Art

Demo

Extra Slides

## AD by Re-execution

- In the sequential context, the tape is elegantly modeled by means of powerful programming abstractions (closures, co-routines), but these are not suitable for GPU execution;
- In our **nested-parallel** context, the tape is **complex/irregular** and must be passed across **deeply nested scopes**.  
Challenging to implement efficiently.

## AD by Re-execution

- In the sequential context, the tape is elegantly modeled by means of powerful programming abstractions (closures, co-routines), but these are not suitable for GPU execution;
- In our **nested-parallel** context, the tape is **complex/irregular** and must be passed across **deeply nested scopes**. Challenging to implement efficiently.

### Key idea #1: Re-execution

Whenever a new scope is entered, the code generation of the return sweep first re-executes the primal trace of that scope.

- **Asymptotics-preserving**: re-execution overhead is constant for non-recursive programs;
- No overhead for **perfectly nested scopes** (other than loops);
- Loops require **checkpointing** & Loop stripmining provides an easy way to navigate an effective **space-time tradeoff**.



# Asymptotic Preserving & No Overhead for Perfect Nests

**Original/Primal Code** is a perfect nest of depth 4:

```
let xss = map ( $\lambda c as \rightarrow$  if  $c$   
                then ...  
                else map ( $\lambda a \rightarrow a*a$ )  $as$   
                )  $cs$   $ass$ 
```



Motivating Example for AD

Gentle Introduction to AD

Key Idea #1: Redundant Execution Instead of Tape

**Key Idea #2: Parallel Constructs are Differentiated at a High Level**

Key Idea #3: Translate Accumulators to Specialized Constructs

Experimental Results: Competitive with State of the Art

Demo

Extra Slides

# Differentiating Reduce at a High Level

- Reduce combines all elements of an array with a binary associative operator  $\odot$ :

$$\mathbf{let} \ y = \mathbf{reduce} \ \odot \ e_{\odot} \ [a_0, a_1, \dots, a_{n-1}]$$

$\equiv$

$$\mathbf{let} \ y = a_0 \odot a_1 \odot \dots \odot a_{n-1}$$

# Differentiating Reduce at a High Level

- Reduce combines all elements of an array with a binary associative operator  $\odot$ :

$$\mathbf{let } y = \mathbf{reduce } \odot e_{\odot} [a_0, a_1, \dots, a_{n-1}]$$

$\equiv$

$$\mathbf{let } y = a_0 \odot a_1 \odot \dots \odot a_{n-1}$$

- For each  $a_i$  in the array, we can group the terms of reduce as:

$$\mathbf{let } y = \underbrace{a_0 \odot \dots \odot a_{i-1}}_{l_i} \odot a_i \odot \underbrace{a_{i+1} \odot \dots \odot a_{n-1}}_{r_i}$$

# Differentiating Reduce at a High Level

- Reduce combines all elements of an array with a binary associative operator  $\odot$ :

$$\mathbf{let} \ y = \mathbf{reduce} \ \odot \ e_{\odot} \ [a_0, a_1, \dots, a_{n-1}]$$

$\equiv$

$$\mathbf{let} \ y = a_0 \odot a_1 \odot \dots \odot a_{n-1}$$

- For each  $a_i$  in the array, we can group the terms of reduce as:

$$\mathbf{let} \ y = \underbrace{a_0 \odot \dots \odot a_{i-1}}_{l_i} \odot a_i \odot \underbrace{a_{i+1} \odot \dots \odot a_{n-1}}_{r_i}$$

And then directly apply the AD rewrite rule

$$\bar{a}_i \ \bar{+} = \frac{\partial(l_i \odot a_i \odot r_i)}{\partial a_i} \ \bar{y}$$

- We compute all  $\bar{a}_i$  in parallel by mapping over all instances of  $l_i, a_i, r_i$ . How do we compute all instances of  $l_i$  and  $r_i$ ?

## Computing $l_i$ and $r_i$

- For each  $i \in \{0, \dots, n-1\}$ , need to compute  $l_i$  and  $r_i$

$$\underbrace{a_0 \odot \dots \odot a_{i-1}}_{l_i} \odot a_i \odot \underbrace{a_{i+1} \odot \dots \odot a_{n-1}}_{r_i}$$

- Compute all  $l_i$  with a (parallel) prefix sum operation, a.k.a., scan:

$$\begin{aligned} \text{let } ls &= \mathbf{scan} \odot e_{\odot} [a_0, a_1, \dots, a_{n-1}] \\ &\equiv \left[ \underbrace{e_{\odot}}_{l_0}, \underbrace{a_0}_{l_1}, \underbrace{a_0 \odot a_1}_{l_2}, \dots, \underbrace{a_0 \odot \dots \odot a_{n-2}}_{l_{n-1}} \right] \end{aligned}$$

- For the  $r_i$ s, do a backward scan, i.e., on the reversed array

$$\begin{aligned} \text{let } rs &= \mathbf{reverse} \ as \triangleright \mathbf{scan} (\lambda x y \rightarrow y \odot x) \ e_{\odot} [a_0, a_1, \dots, a_{n-1}] \\ &\triangleright \mathbf{reverse} \\ &\equiv \left[ \underbrace{a_0 \odot \dots \odot a_{n-2}}_{r_0}, \dots, \underbrace{a_{n-2} \odot a_{n-1}}_{r_{n-3}}, \underbrace{a_{n-1}}_{r_{n-2}}, \underbrace{e_{\odot}}_{r_{n-1}} \right] \end{aligned}$$

# The Reduce Rule

- The differentiation of reduce results in the following statements

**let**  $y = \mathbf{reduce} \odot e_{\odot} [a_0, a_1, \dots, a_{n-1}]$

$\vdots$

**let**  $ls = \mathbf{scan} \odot e_{\odot} as$

**let**  $rs = \mathbf{reverse} as \triangleright \mathbf{scan} (\lambda x y \rightarrow y \odot x) e_{\odot} \triangleright \mathbf{reverse}$

**let**  $\overline{as} \overline{r} = \mathbf{map} \left( \lambda l_i a_i r_i \rightarrow \frac{\partial(l_i \odot a_i \odot r_i)}{\partial a_i} \overline{y} \right) ls as rs$



# The Reduce Rule

- The differentiation of reduce results in the following statements

**let**  $y = \mathbf{reduce} \odot e_{\odot} [a_0, a_1, \dots, a_{n-1}]$

$\vdots$

**let**  $ls = \mathbf{scan} \odot e_{\odot} as$

**let**  $rs = \mathbf{reverse} as \triangleright \mathbf{scan} (\lambda x y \rightarrow y \odot x) e_{\odot} \triangleright \mathbf{reverse}$

**let**  $\overline{as} \overline{r} = \mathbf{map} \left( \lambda l_i a_i r_i \rightarrow \frac{\partial(l_i \odot a_i \odot r_i)}{\partial a_i} \overline{y} \right) ls as rs$

- The rule is asymptotics-preserving: scan has the same asymptotics as reduce.
- Specialized rules for other operators (+, min, max, \*) admit significantly more efficient implementations (map-reduce).

Motivating Example for AD

Gentle Introduction to AD

Key Idea #1: Redundant Execution Instead of Tape

Key Idea #2: Parallel Constructs are Differentiated at a High Level

**Key Idea #3: Translate Accumulators to Specialized Constructs**

Experimental Results: Competitive with State of the Art

Demo

Extra Slides

# Optimizing Accumulators: Matrix-Multiplication Eg

A class of optimizations refers to translating accumulators to more specialized constructs, e.g., reductions.

## Matrix Multiplication Original:

```
forall i = 0..n-1           // map
  forall j = 0 ..n-1       // map
    for k = 0 .. n-1      // map-reduce
      c[i,j] += a[i,k] * b[k,j]
```

# Optimizing Accumulators: Matrix-Multiplication Eg

A class of optimizations refers to translating accumulators to more specialized constructs, e.g., reductions.

## Matrix Multiplication Original:

```
forall i = 0..n-1                // map
  forall j = 0 ..n-1             // map
    for k = 0 .. n-1            // map-reduce
      c[i,j] += a[i,k] * b[k,j]
```

## Matrix Multiplication Reverse-AD with accumulators:

```
forall i = 0..n-1
  forall j = 0 ..n-1
    forall k = 0 .. n-1
       $\bar{a}[i,k] += b[k,j] * \bar{c}[i,j]$ 
       $\bar{b}[k,j] += a[i,k] * \bar{c}[i,j]$ 
```

Rewrite the accumulations as classical reductions. This requires:

# Optimizing Accumulators: Matrix-Multiplication Eg

A class of optimizations refers to translating accumulators to more specialized constructs, e.g., reductions.

## Matrix Multiplication Original:

```
forall i = 0..n-1                // map
  forall j = 0 ..n-1             // map
    for k = 0 .. n-1            // map-reduce
      c[i,j] += a[i,k] * b[k,j]
```

## Matrix Multiplication Reverse-AD with accumulators:

```
forall i = 0..n-1
  forall j = 0 ..n-1
    forall k = 0 .. n-1
       $\bar{a}[i,k] += b[k,j] * \bar{c}[i,j]$ 
       $\bar{b}[k,j] += a[i,k] * \bar{c}[i,j]$ 
```

Rewrite the accumulations as classical reductions. This requires:

- to distribute the loop-nest over each statement
- bring innermost the parallel loop to which the write access is invariant to.

# MMM Example: Optimizing Generalized Reductions

- Distribute the loop-nest over each statement.
- Bring innermost the parallel loop to which the write access is invariant to.

```
forall i = 0..n-1
  forall k = 0 .. n-1
    forall j = 0 ..n-1
       $\bar{a}[i,k] += b[k,j] * \bar{c}[i,j]$ 
```

```
forall k = 0 .. n-1
  forall j = 0 ..n-1
    forall i = 0..n-1
       $\bar{b}[k,j] += a[i,k] * \bar{c}[i,j]$ 
```

**Perform a strength reduction: result can be summed and accumulated once**

# MMM Example: Optimizing Generalized Reductions

- Distribute the loop-nest over each statement.
- Bring innermost the parallel loop to which the write access is invariant to.

```
forall i = 0..n-1
  forall k = 0 .. n-1
    forall j = 0 ..n-1
       $\bar{a}[i,k] += b[k,j] * \bar{c}[i,j]$ 
```

```
forall k = 0 .. n-1
  forall j = 0 ..n-1
    forall i = 0..n-1
       $\bar{b}[k,j] += a[i,k] * \bar{c}[i,j]$ 
```

**Perform a strength reduction: result can be summed and accumulated once**

```
forall i = 0..n-1
  forall k = 0 .. n-1
    acc = 0
    for j = 0 ..n-1
      acc = acc + b[k,j] *  $\bar{c}[i,j]$ 
     $\bar{a}[i,k] += acc$ 
```

```
forall k = 0 .. n-1
  forall j = 0 ..n-1
    acc = 0
    for i = 0..n-1
      acc = acc + a[i,k] *  $\bar{c}[i,j]$ 
     $\bar{b}[k,j] += acc$ 
```

# MMM Example: Optimizing Generalized Reductions

```
forall i = 0..n-1 // map
  forall k = 0 .. n-1 // map
    acc = 0
    for j = 0 ..n-1 // map-reduce
      acc = acc + b[k,j] *  $\bar{c}[i,j]$ 
       $\bar{a}[i,k]$  += acc

forall k = 0 .. n-1 // map
  forall j = 0 ..n-1 // map
    acc = 0
    for i = 0..n-1 // map-reduce
      acc = acc + a[i,k] *  $\bar{c}[i,j]$ 
       $\bar{b}[k,j]$  += acc
```

Now re-write (most of) the accumulations as classical reductions:

```
map(\ i ->
  map(\ k ->
    let acc = map2 (*) b[k]  $\bar{c}[i]$ 
      |> reduce (+) 0
    let  $\bar{a}[i,k]$  += acc in  $\bar{a}$ 
  ) (iota n)
) (iota n)

map(\ k ->
  map(\ j ->
    let acc = map2 (*) a[:,k]  $\bar{c}[:,j]$ 
      |> reduce (+) 0
    let  $\bar{b}[k,j]$  += acc in  $\bar{b}$ 
  ) (iota n)
) (iota n)
```

In this form, the (enhanced) Futhark compiler would apply block and register tiling, and also implement the technique of parallelizing the innermost dimension proposed by [Rasch,Schulze, and Gorlatch, 2019]



Motivating Example for AD

Gentle Introduction to AD

Key Idea #1: Redundant Execution Instead of Tape

Key Idea #2: Parallel Constructs are Differentiated at a High Level

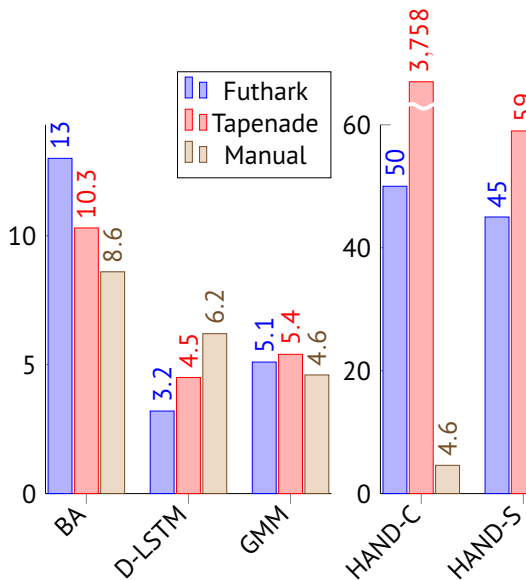
Key Idea #3: Translate Accumulators to Specialized Constructs

**Experimental Results: Competitive with State of the Art**

Demo

Extra Slides

# CPU Benchmarks - ADBench



- ADBench: a collection of AD benchmarks for comparing sequential AD tools.

- Benchmarked Futhark using its C backend.

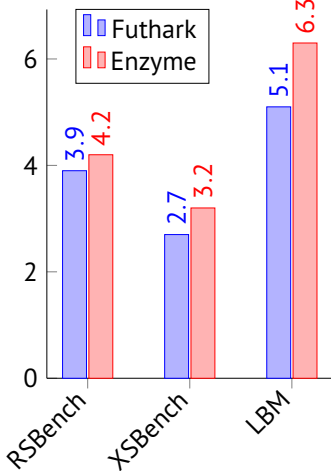
- Performance measured in **AD overhead**:

$$\frac{\text{differentiated runtime}}{\text{original runtime}}$$

**the lower the better**

- On BA, bottleneck is due to packing the result in CSR.

# GPU Benchmarks - vs. Enzyme

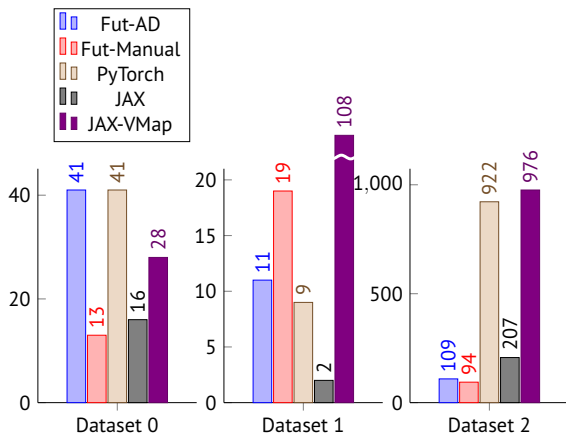


- Performance measured in **AD overhead:**

$$\frac{\text{differentiated runtime}}{\text{original runtime}}$$

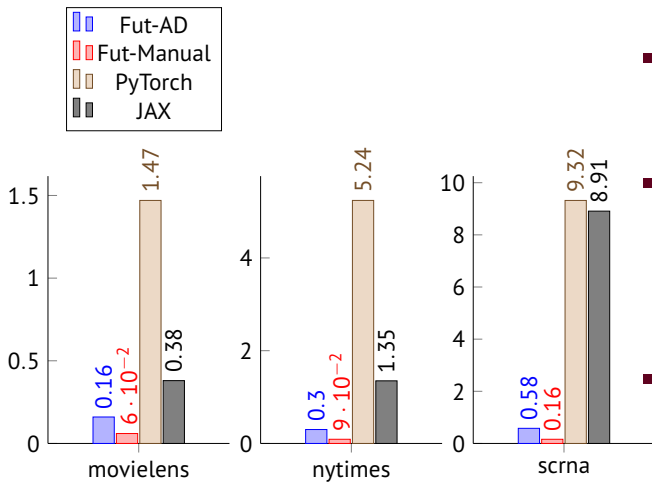
- Enzyme is state-of-the-art LLVM compiler plugin that performs AD on a low-level imperative IR.
- RSbench and XSBench are comprised of a large parallel loop with inner sequential loops and branches.
- LBM consists of a large sequential loop containing a parallel loop.

# *k*-means by Newton's Method on NVIDIA's A100 GPU



- Performance measured in **milliseconds**.
- PyTorch and JAX use hand-tuned matrix primitives;
- JAX(vmap) instead uses JAX's vectorizing map operation.

# Sparse $k$ -means on Nvidia's A100 GPU



- Performance measured in **seconds**.
- PyTorch and JAX use hand-tuned matrix primitives and sparse libraries.
- Futhark just uses a standard CSR implementation.

# GMM: GPU Performance vs PyTorch on A100 & MI100

We test benchmark GMM from ADBench on 32-bit floats.  
This is neutral ground.

	<b>Measurement</b>	<b>D<sub>0</sub></b>	<b>D<sub>1</sub></b>	<b>D<sub>2</sub></b>	<b>D<sub>3</sub></b>	<b>D<sub>4</sub></b>	<b>D<sub>5</sub></b>
A100	<b>PyT. Jacob. (ms)</b>	7.4	15.8	15.2	5.9	12.5	64.8
	<b>Fut. Speedup</b>	2.1	2.2	1.4	1.6	1.5	1.0
	<b>PyT. Overhead</b>	3.5	4.9	2.8	3.2	4.0	3.2
	<b>Fut. Overhead</b>	2.0	1.8	1.9	2.7	2.8	2.8
MI100	<b>PyT. Jacob. (ms)</b>	20.9	51.5	42.5	20.7	38.5	193.1
	<b>Fut. Speedup</b>	3.3	4.0	2.1	2.9	2.5	1.7
	<b>PyT. Overhead</b>	5.9	5.3	2.4	2.6	3.1	2.8
	<b>Fut. Overhead</b>	3.0	2.9	3.0	2.8	2.8	2.8

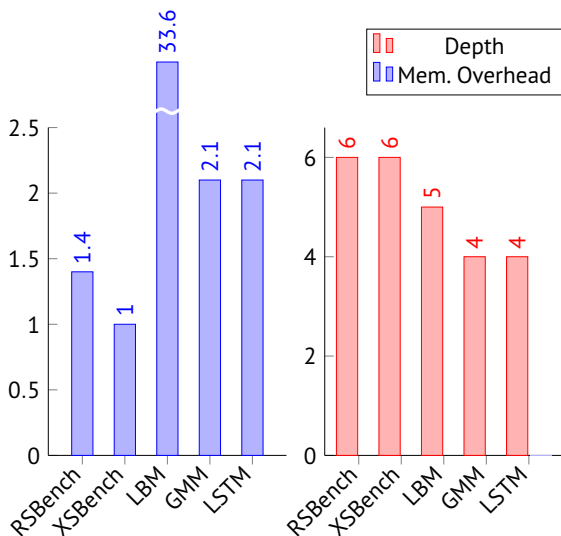
# LSTM: GPU Performance vs PyTorch & JAX

This is PyTorch's home ground, because matrix-multiplications take about 75% of runtime, and matrix-multiplication is a primitive in PyTorch, while Futhark needs to work hard for it.

		Speedups				
		PyTorch Jacob.	Futhark	nn.LSTM	JAX	JAX(vmap)
A100	D <sub>0</sub>	45.4 ms	3.0	11.6	4.5	0.3
	D <sub>1</sub>	740.1 ms	3.3	22.1	6.4	0.9
MI100	D <sub>0</sub>	89.8 ms	2.6	4.0	—	—
	D <sub>1</sub>	1446.9 ms	1.8	5.4	—	—
		Overheads				
		PyTorch	Futhark	nn.LSTM	JAX	JAX(vmap)
A100	D <sub>0</sub>	4.1	2.1	2.7	3.5	1.4
	D <sub>1</sub>	4.3	3.9	2.2	3.7	0.8
MI100	D <sub>0</sub>	5.0	4.2	7.2	—	—
	D <sub>1</sub>	7.9	3.9	6.6	—	—

cuDNN-based refers to the (manual) `torch.nn.LSTM` library.

# GPU Benchmarks - Depth and Memory Consumption



■ **AD Memory overhead:**

$$\frac{\text{differentiated mem.}}{\text{original mem.}}$$

■ **Loop strip-mining** reduces LBM's memory overhead to  $8.7\times$ , at the cost of  $1.3\times$  increase in runtime.

**Strong performance on programs with non-trivial depth demonstrates the viability of a recomputation-based approach to AD.**



DEMO

# Reduce with Multiplication

$$\text{let } y = w_0 \cdot w_1 \cdot \dots \cdot w_{n-1}$$

## Sequential:

```
loop y = 1
  for i = 0 .. n-1 do
    y * ws[i]
```

## Parallel:

```
let y = reduce (*) 1 ws
```

## Why Reduce with Arbitrary Operators?

$$\underbrace{\sum_{i,j=1..k}^{j<i} w_j \cdot w_i}_{P_k} = \underbrace{\left(\sum_{i,j=1..k-1}^{j<i} w_j \cdot w_i\right)}_{P_{k-1}} + \underbrace{\left(\sum_{i=1..k-1} w_i\right)}_{S_{k-1}} \cdot w_k$$

# Why Reduce with Arbitrary Operators?

$$\underbrace{\sum_{i,j=1..k}^{j<i} w_j \cdot w_i}_{P_k} = \underbrace{\left(\sum_{i,j=1..k-1}^{j<i} w_j \cdot w_i\right)}_{P_{k-1}} + \underbrace{\left(\sum_{i=1..k-1} w_i\right)}_{S_{k-1}} \cdot w_k$$

## Sequential:

```
let seqop (P,S) w_new = (P + S*w_new, S + w_new)
```

```
...
```

```
let (res, _) =  
  loop y = (0,0)  
    for i = 0 .. n-1 do  
      seqop y ws[i]
```

# Why Reduce with Arbitrary Operators?

$$\underbrace{\sum_{i,j=1..k}^{j<i} w_j \cdot w_i}_{P_k} = \underbrace{\left(\sum_{i,j=1..k-1}^{j<i} w_j \cdot w_i\right)}_{P_{k-1}} + \underbrace{\left(\sum_{i=1..k-1} w_i\right)}_{S_{k-1}} \cdot w_k$$

## Sequential:

```
let seqop (P,S) w_new = (P + S*w_new, S + w_new)
```

...

```
let (res, _) =  
  loop y = (0,0)  
    for i = 0 .. n-1 do  
      seqop y ws[i]
```

## Parallel:

```
let parop (P1, S1) (P2, S2) = (P1 + P2 + S1*S2, S1+S2)
```

```
let (res, _) = map (\w -> (0, w)) ws — lifting  
  |> reduce parop (0,0)
```

# Conclusions

- AD in a **nested-parallel, high-level** and **hardware-neutral** functional language.
- **Key idea:** high-level differentiation using specialized rules for parallel combinators.
- **Key idea:** re-computation instead of a tape (except for loops!).
- Strong performance against state-of-the-art AD competitors.
- The implementation is available now in the Futhark compiler—try it out!

<https://futhark-lang.org>

Motivating Example for AD

Gentle Introduction to AD

Key Idea #1: Redundant Execution Instead of Tape

Key Idea #2: Parallel Constructs are Differentiated at a High Level

Key Idea #3: Translate Accumulators to Specialized Constructs

Experimental Results: Competitive with State of the Art

Demo

**Extra Slides**

# Loops

- **Sequential** loops are sugar for tail-recursive functions.
- **Loop parameters** are **variables** which are variant through the loop and are returned as the result of the loop.

```
loop y = 2 for  $i = 0 \dots n - 1$  do  
  let  $y' = y * y$   
  in  $y'$ 
```

```
y = 2  
for  $i = 0 \dots n - 1$  do  
   $y = y * y$ 
```

(Imperative analog)

- Storing the loop parameter **y** on the tape for each iteration is required in order to be able to execute the loop backwards.
- Parallel constructs do not require such check-pointing.



# Differentiating Loops

```
let  $y'' =$   
  loop  $y = y_0$  for  $i = 0 \dots n - 1$  do  
     $stms_{loop}$   
  in  $y'$ 
```

1. Re-execute the original loop, save the value of  $y$  in each iteration in  $ys$ .

```
2 let  $ys_0 = scratch(n,$   
3            $sizeof(y_0))$   
4 let  $(y'', ys) =$   
5 loop  $(y, ys) = (y_0, ys_0)$   
6 for  $i = 0 \dots n - 1$  do  
7   let  $ys[i] = y$   
8    $stms_{loop}$   
9   in  $(y', ys)$   
12 let  $(\overline{y''''}, \overline{fvs_l}) =$   
13 loop  $(\overline{y}, \overline{fvs_l}) = (\overline{y''}, \overline{fvs_{l_0}})$   
14 for  $i = n - 1 \dots 0$  do  
15   let  $y = ys[i]$   
16    $\overrightarrow{stms_{loop}}$   
17    $\overleftarrow{stms_{loop}}$   
18   in  $(\overline{y'}, \overline{fvs'_l})$   
19 let  $\overline{y_0} += \overline{y''''}$ 
```

Primal Trace

Return Sweep

# Differentiating Loops

```
let  $y'' =$   
  loop  $y = y_0$  for  $i = 0 \dots n - 1$  do  
     $stms_{loop}$   
  in  $y'$ 
```

1. Re-execute the original loop, save the value of  $y$  in each iteration in  $ys$ .
2. Compute the adjoint contributions of the loop.

```
2 let  $ys_0 = \text{scratch}(n,$   
3      $\text{sizeof}(y_0))$   
4 let  $(y'', ys) =$   
5 loop  $(y, ys) = (y_0, ys_0)$   
6 for  $i = 0 \dots n - 1$  do  
7   let  $ys[i] = y$   
8    $stms_{loop}$   
9   in  $(y', ys)$   
12 let  $(\overline{y'''}, \overline{fvs_l}) =$   
13 loop  $(\overline{y}, \overline{fvs_l}) = (\overline{y''}, \overline{fvs_{l_0}})$   
14 for  $i = n - 1 \dots 0$  do  
15   let  $y = ys[i]$   
16    $\overrightarrow{stms_{loop}}$   
17    $\overleftarrow{stms_{loop}}$   
18   in  $(\overline{y'}, \overline{fvs'_l})$   
19 let  $\overline{y_0} += \overline{y'''}$ 
```

Primal Trace

Return Sweep

# Differentiating Loops

```
let  $y'' =$   
  loop  $y = y_0$  for  $i = 0 \dots n - 1$  do  
     $stms_{loop}$   
  in  $y'$ 
```

1. Re-execute the original loop, save the value of  $y$  in each iteration in  $ys$ .
2. Compute the adjoint contributions of the loop.
  - ▶ Run the loop backwards

```
2 let  $ys_0 = \text{scratch}(n,$   
3      $\text{sizeof}(y_0))$   
4 let  $(y'', ys) =$   
5 loop  $(y, ys) = (y_0, ys_0)$   
6 for  $i = 0 \dots n - 1$  do  
7   let  $ys[i] = y$   
8    $stms_{loop}$   
9   in  $(y', ys)$   
12 let  $(\overline{y'''}, \overline{fvs_l}) =$   
13 loop  $(\overline{y}, \overline{fvs_l}) = (\overline{y''}, \overline{fvs_{l_0}})$   
14 for  $i = n - 1 \dots 0$  do  
15   let  $y = ys[i]$   
16    $\overrightarrow{stms_{loop}}$   
17    $\overleftarrow{stms_{loop}}$   
18   in  $(\overline{y'}, \overline{fvs'_l})$   
19 let  $\overline{y_0} += \overline{y'''}$ 
```

Primal Trace

Return Sweep

# Differentiating Loops

```
let  $y'' =$   
  loop  $y = y_0$  for  $i = 0 \dots n - 1$  do  
     $stms_{loop}$   
  in  $y'$ 
```

1. Re-execute the original loop, save the value of  $y$  in each iteration in  $ys$ .
2. Compute the adjoint contributions of the loop.
  - ▶ Run the loop backwards
  - ▶ Restore the value of  $y$  from  $ys$

```
2 let  $ys_0 = \text{scratch}(n,$   
3      $\text{sizeof}(y_0))$   
4 let  $(y'', ys) =$   
5 loop  $(y, ys) = (y_0, ys_0)$   
6 for  $i = 0 \dots n - 1$  do  
7   let  $ys[i] = y$   
8    $stms_{loop}$   
9   in  $(y', ys)$   
12 let  $(\overline{y'''}, \overline{fvs_l}) =$   
13 loop  $(\overline{y}, \overline{fvs_l}) = (\overline{y''}, \overline{fvs_{l_0}})$   
14 for  $i = n - 1 \dots 0$  do  
15   let  $y = ys[i]$   
16    $\overrightarrow{stms_{loop}}$   
17    $\overleftarrow{stms_{loop}}$   
18   in  $(\overline{y'}, \overline{fvs'_l})$   
19 let  $\overline{y_0} += \overline{y'''}$ 
```

Primal Trace

Return Sweep

# Differentiating Loops

```
let  $y'' =$   
  loop  $y = y_0$  for  $i = 0 \dots n - 1$  do  
     $stms_{loop}$   
  in  $y'$ 
```

1. Re-execute the original loop, save the value of  $y$  in each iteration in  $ys$ .
2. Compute the adjoint contributions of the loop.
  - ▶ Run the loop backwards
  - ▶ Restore the value of  $y$  from  $ys$
  - ▶ Re-execute the body of the original loop

```
2 let  $ys_0 = \text{scratch}(n,$   
3      $\text{sizeof}(y_0))$   
4 let  $(y'', ys) =$   
5 loop  $(y, ys) = (y_0, ys_0)$   
6 for  $i = 0 \dots n - 1$  do  
7   let  $ys[i] = y$   
8    $stms_{loop}$   
9   in  $(y', ys)$   
12 let  $(\overline{y'''}, \overline{fvs_l}) =$   
13 loop  $(\overline{y}, \overline{fvs_l}) = (\overline{y''}, \overline{fvs_{l_0}})$   
14 for  $i = n - 1 \dots 0$  do  
15   let  $y = ys[i]$   
16    $\overrightarrow{stms_{loop}}$   
17    $\overleftarrow{stms_{loop}}$   
18   in  $(\overline{y'}, \overline{fvs'_l})$   
19 let  $\overline{y_0} += \overline{y'''}$ 
```

Primal Trace

Return Sweep

# Differentiating Loops

```
let  $y'' =$   
  loop  $y = y_0$  for  $i = 0 \dots n - 1$  do  
     $stms_{loop}$   
  in  $y'$ 
```

1. Re-execute the original loop, save the value of  $y$  in each iteration in  $ys$ .
2. Compute the adjoint contributions of the loop.
  - ▶ Run the loop backwards
  - ▶ Restore the value of  $y$  from  $ys$
  - ▶ Re-execute the body of the original loop
  - ▶ Compute the adjoints of the body

```
2 let  $ys_0 = \text{scratch}(n,$   
3      $\text{sizeof}(y_0))$   
4 let  $(y'', ys) =$   
5 loop  $(y, ys) = (y_0, ys_0)$   
6 for  $i = 0 \dots n - 1$  do  
7   let  $ys[i] = y$   
8    $stms_{loop}$   
9   in  $(y', ys)$   
12 let  $(\overline{y'''}, \overline{fvs_l}) =$   
13 loop  $(\overline{y}, \overline{fvs_l}) = (\overline{y''}, \overline{fvs_{l_0}})$   
14 for  $i = n - 1 \dots 0$  do  
15   let  $y = ys[i]$   
16    $\overrightarrow{stms_{loop}}$   
17    $\overleftarrow{stms_{loop}}$   
18   in  $(\overline{y'}, \overline{fvs'_l})$   
19 let  $\overline{y_0} += \overline{y'''}$ 
```

Primal Trace

Return Sweep

# Loop Strip-mining implements the Space-Time Trade-off

- **Loop strip-mining** partitions a loop into a loop nest

```
loop  $y = y_0$   
  for  $i = 0 \dots n^3 - 1$  do  
    stms  
⇒  
  loop  $y_j = y_0$  for  $j = 0 \dots n - 1$  do  
    loop  $y_k = y_j$  for  $k = 0 \dots n - 1$  do  
      loop  $y_m = y_k$  for  $m = 0 \dots n - 1$  do  
        let  $i = j * n^2 + k * n + m$   
        stms
```

- For the original loop, we save  $n^3$  versions of  $y$  on the tape.
- For the strip-mined loop, only  $3n$  versions are saved.
- Strip-mining is controlled by the user via a simple annotation.

Strip-mining a loop  $k$  times results in **up to  $k \times$  slowdown**, but memory overhead **decreases** from a factor of  $n^k \times$  to  $n \cdot k \times$ .

## Summary: Tape vs Redundant Execution

- whenever a new scope is entered, the code generation of the reverse trace first re-executes the primal trace of that scope;
- the re-execution overhead is at worst proportional with the deepest scope nest of the program (which is constant);
- “the tape” is part of the program and subject to aggressive optimizations (especially in a purely functional context);
- in most cases, it is more efficient to re-compute scalars rather than to access them from the tape (global memory);
- **loops require checkpointing, parallel constructs do not.**
- Subject to checkpointing are only the loops appearing directly in the current scope of the reverse-trace code generation (i.e., inner loops of the primal trace are not).



# Map without Free Variables

- Map is equivalent to a parallel for-loop

**let**  $xs = \mathbf{map}$   $(\lambda a\ b \rightarrow \mathbf{let}\ res = a * b\ \mathbf{in}\ res)$   $as\ bs$

$\Updownarrow$

**forall**  $i = 0 \dots n - 1$

$xs[i] = as[i] * bs[i]$

# Map without Free Variables

- Map is equivalent to a parallel for-loop

```
let xs = map ( $\lambda a b \rightarrow$  let res =  $a * b$  in res) as bs
```

$\Updownarrow$

```
forall  $i = 0 \dots n - 1$ 
```

```
   $xs[i] = as[i] * bs[i]$ 
```

- Differentiating map is straightforward in the absence of free variables

```
let  $\overline{as}, \overline{bs} = \mathbf{map} (\lambda a b \overline{x} \overline{a_0} \overline{b_0} \rightarrow$ 
```

```
  let res =  $a * b$ 
```

```
  let  $\overline{a} = \overline{a_0} + b * \overline{x}$ 
```

```
  let  $\overline{b} = \overline{b_0} + a * \overline{x}$ 
```

```
  in  $\overline{a}, \overline{b})$  as bs  $\overline{xs} \overline{as_0} \overline{bs_0}$ 
```

# Map with Free Variables

- Maps involving **free variables** are more complicated to differentiate

**let**  $xs = \text{map } (\lambda a \rightarrow a * b) \text{ as}$

- Naive approach: turn free variables into bound variables.

**let**  $xs = \text{map } (\lambda a \ b' \rightarrow a * b') \text{ as } (\text{replicate } n \ b)$

- Problem: asymptotically inefficient for partially used free arrays.

# Efficient Maps with Free Variables

- In an impure language, asymptotics-preserving adjoint updates for free array variables can be implemented as a **generalized reduction**:
  - ▶ preserves the useful properties of maps (parallel loops),
  - ▶ generalization of map, reduce, reduce-by-index, scatter.
- In this setting, the adjoint of a free array variable  $as[i]$  can be updated with an operation  $\overline{as}[i] += v$ .
- In our pure setting, we introduce **accumulators**.
  - ▶ **Write-only** view of an array.
  - ▶ Guarantees the generalized reduction properties at the type level.
- **An important set of optimizations refer to specializing generalized reductions to maps, reduce, reduce-by-index.**

# Handling of Loops with In-Place Updates

The brownian loop bridge from OptionPricing, FinPar suite:

```
let bbrow =  
  loop bbrow for i in 1..<num_dates do  
    let bbrow[ bi[i] - 1 ] =  
      sd[ i ] * gauss[ i ] +  
      rw[ i ] * bbrow[ ri[i] - 1 ] +  
      lw[ i ] * bbrow[ li[i] - 1 ]  
  in bbrow
```

**Checkpointing `bbrow` for each iteration breaks work asymptotics!**

**If the loop exhibits only true (RAW) dependencies, then:**

# Handling of Loops with In-Place Updates

The brownian loop bridge from OptionPricing, FinPar suite:

```
let bbrow =  
  loop bbrow for i in 1..<num_dates do  
    let bbrow[ bi[i] - 1 ] =  
      sd[ i ] * gauss[ i ] +  
      rw[ i ] * bbrow[ ri[i] - 1 ] +  
      lw[ i ] * bbrow[ li[i] - 1 ]  
  in bbrow
```

**Checkpointing `bbrow` for each iteration breaks work asymptotics!**

**If the loop exhibits only true (RAW) dependencies, then:**

- On primal: before the loop, checkpoint the indices of `bbrow` written through the loop (e.g., the whole array);
- On return: after computing the adjoint of the loop, restore the `bbrow` array to the state before the loop.

**Rationale:**

# Handling of Loops with In-Place Updates

## The brownian loop bridge from OptionPricing, FinPar suite:

```
let bbrow =  
  loop bbrow for i in 1..<num_dates do  
    let bbrow[ bi[i] - 1 ] =  
      sd[ i ] * gauss[ i ] +  
      rw[ i ] * bbrow[ ri[i] - 1 ] +  
      lw[ i ] * bbrow[ li[i] - 1 ]  
  in bbrow
```

**Checkpointing `bbrow` for each iteration breaks work asymptotics!**

**If the loop exhibits only true (RAW) dependencies, then:**

- On primal:** before the loop, checkpoint the indices of `bbrow` written through the loop (e.g., the whole array);
- On return:** after computing the adjoint of the loop, restore the `bbrow` array to the state before the loop.

**Rationale:**

- Parallel prog means avoiding (cross-iteration) dependencies;
- WAR & WAW are named **false** dependencies for a reason!