# Reverse-Mode Automatic Differentiation in Futhark Language

#### Cosmin E. Oancea (cosmin.oancea@diku.dk) work in collaboration with Troels Henriksen, Ola Rønning and Robert Schenck

Department of Computer Science (DIKU) University of Copenhagen

24th of July 2023, MODE Workshop

# **Overview**



- This talk presents an AD algorithm [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.
- AD is applied **before** parallelism is mapped to the hardware.

Work is aimed to explore how the richer semantics of high-level parallel language enables AD as a first-lass language citizen.

[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. Reverse-AD: Core Re-Write Rule

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

Extra Slides Motivating Example for AD Differentiating Loops & Time-Space Tradeoff Differentiating Map Loops with In-Place Updates

#### Reverse-AD: Core Re-Write Rule

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

Extra Slides Motivating Example for AD Differentiating Loops & Time-Space Tradeoff Differentiating Map Loops with In-Place Updates

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

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

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

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

Initially  $\overline{y} = \frac{\partial y}{\partial y} = 1$ , and eventually the adjoints of the input  $\overline{x_i} = \frac{\partial y}{\partial x_i}$  are computed by applying the following re-write rule:

Primal
$$\bigvee_{v}$$
Reverse $\widehat{x} = \overline{x} + \frac{\partial f(x, y, ...)}{\partial x} \cdot \overline{z}$ Trace: $\downarrow_{v}$ Sweep: $\cdots$  $\bigvee_{v}$  $\cdots$  $\cdots$  $\cdots$ 

(\*) final value of  $\overline{z}$  is known;

(\*\*)  $\overline{z}$  is dead after this point;  $\overline{x}$  and  $\overline{y}$  are still under computation.

**A tape** is needed to recover the values of x and y (used in  $\frac{\partial f(x,y)}{\partial x}$ )

Adjoint of var t is  $\overline{t} \equiv \frac{\partial y}{\partial t}$ . Goal: compute the adjoints of the input.Core Rewrite Rule:  $z = f(x, y) \Rightarrow$  $\overline{x} + = \frac{\partial f(x, y)}{\partial x} \cdot \overline{z}$  $\overline{y} + = \frac{\partial f(x, y)}{\partial x} \cdot \overline{z}$  $\overline{y} + = \frac{\partial f(x, y)}{\partial x} \cdot \overline{z}$  $P(x_0, x_1):$  $t_0 = \sin(x_0)$  $t_0 = \sin(x_0)$  $t_1 = x_1 \cdot t_0$  $y = x_0 + t_1$  $y = x_0 + t_1$ return y $y = x_0 + t_1$ 

**Adjoint** of var *t* is  $\overline{t} \equiv \frac{\partial y}{\partial t}$ . **Goal:** compute the adjoints of the input. **Core Rewrite Rule**:  $y = f(x_1, x_2) \Rightarrow \frac{\overline{x_1} + = \frac{\partial f(x_1, x_2)}{\partial x_1} \cdot \overline{z}}{\overline{x_2} + = \frac{\partial f(x_1, x_2)}{\partial x_2} \cdot \overline{z}}$ 

$$P(x_0, x_1):$$

$$t_0 = \sin(x_0)$$

$$t_1 = x_1 \cdot t_0$$

$$y = x_0 + t_1$$
return y
$$P'(x_0, x_1):$$

$$t_0 = \sin(x_0)$$

$$t_1 = x_1 \cdot t_0$$

$$y = x_0 + t_1$$

$$\overline{y} = 1; \quad \overline{x_0} = 0; \quad \overline{t_1} = 0$$

$$\overline{x_0} = \overline{x_0} + \frac{\partial(x_0 + t_1)}{\partial t_0} \cdot \overline{y} = 1 \cdot 1 = 1$$

$$\overline{t_0} = \overline{t_0} + \frac{\partial(x_0 + t_1)}{\partial t_1} \cdot \overline{y} = 1 \cdot 1 = 1$$

**Adjoint** of var *t* is  $\overline{t} \equiv \frac{\partial y}{\partial t}$ . **Goal:** compute the adjoints of the input. **Core Rewrite Rule**:  $y = f(x_1, x_2) \Rightarrow \frac{\overline{x_1} + = \frac{\partial f(x_1, x_2)}{\partial x_1} \cdot \overline{z}}{\overline{x_2} + = \frac{\partial f(x_1, x_2)}{\partial x_2} \cdot \overline{z}}$ 

**n**//

$$P(x_{0}, x_{1}):$$

$$t_{0} = \sin(x_{0})$$

$$t_{1} = x_{1} \cdot t_{0}$$

$$y = x_{0} + t_{1}$$
return y
$$P'(x_{0}, x_{1}):$$

$$t_{0} = \sin(x_{0})$$

$$t_{1} = x_{1} \cdot t_{0}$$

$$y = x_{0} + t_{1}$$

$$\overline{y} = 1; \quad \overline{x_{0}} = 0; \quad \overline{t_{1}} = 0$$

$$\overline{x_{0}} = \overline{x_{0}} + \frac{\partial(x_{0} + t_{1})}{\partial x_{0}} \cdot \overline{y} = 1 \cdot 1 = 1$$

$$\overline{t_{1}} = \overline{t_{0}} + \frac{\partial(x_{0} + t_{1})}{\partial t_{1}} \cdot \overline{y} = 1 \cdot 1 = 1$$

$$\overline{x_{1}} = 0; \quad \overline{t_{0}} = 0$$

$$\overline{x_{1}} = \overline{x_{1}} + \frac{\partial(x_{1} \cdot t_{0})}{\partial t_{1}} \cdot \overline{t_{1}} = \overline{x_{1}} + t_{0} \cdot \overline{t_{1}}$$

$$\overline{t_{0}} = \overline{t_{0}} + \frac{\partial(x_{1} \cdot t_{0})}{\partial t_{0}} \cdot \overline{t_{1}} = \overline{t_{0}} + x_{1} \cdot \overline{t_{1}}$$

**Adjoint** of var *t* is  $\bar{t} \equiv \frac{\partial y}{\partial t}$ . **Goal:** compute the adjoints of the input. Core Rewrite Rule:  $y = f(x_1, x_2) \Rightarrow \frac{\overline{x_1} + = \frac{\partial f(x_1, x_2)}{\partial x_1} \cdot \overline{z}}{\overline{x_2} + = \frac{\partial f(x_1, x_2)}{\partial x_2} \cdot \overline{z}}$  $P'(x_0, x_1)$ :  $t_0 = \sin(x_0)$  $P(x_0, x_1)$ :  $t_1 = x_1 \cdot t_0$  $t_0 = \sin(x_0)$  $\implies$  $t_1 = x_1 \cdot t_0$  $v = x_0 + t_1$  $\overline{v} = 1$ :  $\overline{x_0} = 0$ :  $\overline{t_1} = 0$  $y = x_0 + t_1$  $\overline{x_0} = \overline{x_0} + \frac{\partial(x_0 + t_1)}{\partial x_0} \cdot \overline{y} = 1 \cdot 1 = 1$ return y  $\overline{t_1} = \overline{t_0} + \frac{\partial (x_0 + t_1)}{\partial t_1} \cdot \overline{y} = 1 \cdot 1 = 1$  $\overline{x_1} = 0; \quad \overline{t_0} = 0$  $\overline{x_1} = \overline{x_1} + \frac{\partial(x_1 \cdot t_0)}{\partial x_1} \cdot \overline{t_1} = \overline{x_1} + t_0 \cdot \overline{t_1}$  $\overline{t_0} = \overline{t_0} + \frac{\partial (x_1 \cdot \overline{t_0})}{\partial t_0} \cdot \overline{t_1} = \overline{t_0} + x_1 \cdot \overline{t_1}$  $\overline{x_0} = \overline{x_0} + \frac{\partial \sin(x_0)}{\partial x_0} \cdot \overline{t_0} = \overline{x_0} + \cos(x_0) \cdot \overline{t_0}$ return  $\overline{x_0}$ ,  $\overline{x_1}$ 

# **Classical API for AD**

**Reverse-mode:** 

$$\mathsf{vjp}: (f: \alpha \to \beta) \to (\mathbf{x}: \alpha) \to (\overline{\mathbf{y}}: \beta) \to \alpha$$

Forward mode:

$$\mathsf{jvp}: (f: \alpha \to \beta) \to (x: \alpha) \to (dx: \alpha) \to \beta$$

# **Classical API for AD**

**Reverse-mode:** 

$$\mathsf{vjp}: (f: \alpha \to \beta) \to (x: \alpha) \to (\overline{y}: \beta) \to \alpha$$

Forward mode:

$$\mathsf{jvp}: (f: \alpha \to \beta) \to (x: \alpha) \to (dx: \alpha) \to \beta$$

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

#### Reverse-AD: Core Re-Write Rule

#### 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

Extra Slides Motivating Example for AD Differentiating Loops & Time-Space Tradeoff Differentiating Map Loops with In-Place Updates

# The Tape is Challenging to Implement

- In the sequential context, the tape is elegantly modeled by means of powerful programming abstractions (closures, delimited continuations) but these are not suitable for GPU execution;
- In a nested parallel context:
  - Dex: maintains complex and irregular structures of arrays that need to be passed across deeply-nested scopes.
- Enzyme: applies AD to GPU kernels and take advantage of GPU specialized memories but may have problems if kernel's resources are capped out.

#### Implementing the Tape by a Re-Execution Policy

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

- Asymptotics-preserving: re-execution overhead is constant for non-recursive programs, i.e., equal to the depth of the deepest scope nest;
- 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;
- Loop checkpointing can be further optimized by exploiting specific data-dependency properties of loops.

#### Recall: 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:  $y = f(x_1, x_2) \Rightarrow \frac{\overline{x_1} + = \frac{\partial f(x_1, x_2)}{\partial x_1} \cdot \overline{z}}{\overline{x_2} + = \frac{\partial f(x_1, x_2)}{\partial x_2} \cdot \overline{z}}$  $P'(x_0, x_1)$ :  $t_0 = \sin(x_0)$  $P(x_0, x_1)$ :  $t_0 = \sin(x_0)$  $t_1 = x_1 \cdot t_0$  $\implies$  $t_1 = x_1 \cdot t_0$  $v = x_0 + t_1$  $\overline{v} = 1$ :  $\overline{x_0} = 0$ :  $\overline{t_1} = 0$  $y = x_0 + t_1$  $\overline{x_0} = \overline{x_0} + \frac{\partial (x_0 + t_1)}{\partial x_0} \cdot \overline{y} = 1 \cdot 1 = 1$ return y  $\overline{t_1} = \overline{t_0} + \frac{\partial (x_0 + t_1)}{\partial t_1} \cdot \overline{y} = 1 \cdot 1 = 1$  $\overline{x_1} = 0; \quad \overline{t_0} = 0$  $\overline{x_1} = \overline{x_1} + \frac{\partial(x_1 \cdot t_0)}{\partial x_1} \cdot \overline{t_1} = \overline{x_1} + t_0 \cdot \overline{t_1}$  $\overline{t_0} = \overline{t_0} + \frac{\partial (x_1 \cdot \overline{t_0})}{\partial t_0} \cdot \overline{t_1} = \overline{t_0} + x_1 \cdot \overline{t_1}$  $\overline{x_0} = \overline{x_0} + \frac{\partial \sin(x_0)}{\partial x_0} \cdot \overline{t_0} = \overline{x_0} + \cos(x_0) \cdot \overline{t_0}$ 

return  $\overline{x_0}$ ,  $\overline{x_1}$ 

## 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
```

# Asymptotic Preserving & No Overhead for Perfect Nests

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

) ass <del>xss</del>

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

Differentiated Code displays in red the re-execution of the primal:

Reverse-AD: Core Re-Write Rule

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

Extra Slides Motivating Example for AD Differentiating Loops & Time-Space Tradeoff Differentiating Map Loops with In-Place Updates

# **Differentiating Reduce at a High Level**

■ Reduce combines all elements of an array with a binary associative operator ⊙:

$$\begin{array}{rcl} \texttt{let } y \ = \ \texttt{reduce} \ \odot \ e_{\odot} \ [a_0, a_1, \dots, a_{n-1}] \\ & \equiv \\ \texttt{let } y \ = \ a_0 \odot a_1 \odot \dots \odot a_{n-1} \end{array}$$

### Differentiating Reduce at a High Level

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

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

For each *a<sub>i</sub>* in the array, we can group the terms of reduce as:

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

## Differentiating Reduce at a High Level

■ Reduce combines all elements of an array with a binary associative operator ⊙:

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

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

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

And then directly apply the AD rewrite rule

$$\overline{a_i} = \frac{\partial(l_i \odot a_i \odot r_i)}{\partial a_i} \overline{y}$$

We compute all ai in parallel by mapping over all instances of li, ai, ri. How do we compute all instances of li and ri?

## **Computing** *l<sub>i</sub>* and *r<sub>i</sub>*

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

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

Compute all l<sub>i</sub> with a (parallel) prefix sum operation, a.k.a., scan:

$$\begin{array}{rcl} \texttt{let} \ ls = \ \texttt{scan} \ \odot \ e_{\odot} \ [a_0, a_1, \dots, a_{n-1}] \\ \equiv \ [\underbrace{e_{\odot}}_{l_0}, \ \underbrace{a_0}_{l_1}, \ \underbrace{a_0 \odot a_1}_{l_2}, \ \dots, \ \underbrace{a_0 \odot \dots \odot a_{n-2}}_{l_{n-1}}] \end{array}$$

For the *r<sub>i</sub>*s, do a backward scan, i.e., on the reversed array

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

## The Reduce Rule For an Arbitrary Associative Operator

The differentiation of reduce results in the following statements

- The rule is asymptotics-preserving:
  - scan has the same work-depth asymptotics as reduce;
  - requires about 8× more accesses to global memory.

## **Reduce Rule For Commutative and Invertible Operators**

Assume  $\odot : \alpha \to \alpha \to \alpha$  associative and commutative operator. Assume  $\odot$  is also invertible, i.e.,

$$\exists \odot^{-1}$$
 such that  $z \odot^{-1} a = b$  whenever  $b \odot a = z$ 

The differentiation of reduce consists of:

$$\begin{array}{c|c} \textbf{let } y = \textbf{reduce} & \odot & e_{\odot} \left[ a_{0}, a_{1}, \dots, a_{n-1} \right] \\ \vdots \\ \textbf{let } \overline{as} &= \textbf{map}(\lambda \ a_{i} \ \rightarrow \ \textbf{let } b = y \ \odot^{-1} \ a_{i} \\ & \textbf{in} \ \frac{\partial(b \odot a_{i})}{\partial a_{i}} \ \overline{y} \ ) \ as \end{array} \right\}$$
Reverse

 Specialized rules for other operators (+, min, max, \*) admit similar efficient implementations (map-reduce). Reverse-AD: Core Re-Write Rule

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

Extra Slides Motivating Example for AD Differentiating Loops & Time-Space Tradeoff Differentiating Map Loops with In-Place Updates

# **Optimizing Accumulators: Matrix-Multiplication Eg**

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

Matrix Multiplication Original:

# **Optimizing Accumulators: Matrix-Multiplication Eg**

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

Matrix Multiplication Original:

#### Matrix Multiplication Reverse-AD with accumulators:

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:

#### 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.

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.

# 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] * c̄[i,j]
ā[i,k] += acc
```

## MMM Example: Optimizing Generalized Reductions

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

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] Reverse-AD: Core Re-Write Rule

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

Extra Slides Motivating Example for AD Differentiating Loops & Time-Space Tradeoff Differentiating Map Loops with In-Place Updates

## Sequential CPU Benchmarks - ADBench



# GPU Benchmarks - vs. Enzyme



 Performance measured in AD overhead:

differentiated runtime 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 parallell 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 miliseconds.
- PyTorch and JAX use hand-tuned matrix primitives;
- JAX(vmap) instead uses JAX's vectorizing map operation.


- 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.

	Measurement	<b>D</b> <sub>0</sub>	$\mathbf{D}_1$	<b>D</b> <sub>2</sub>	<b>D</b> <sub>3</sub>	$D_4$	<b>D</b> <sub>5</sub>
A100	PyT. Jacob. (ms)	7.4	15.8	15.2	5.9	12.5	64.8
	Fut. Speedup	2.1	2.2	1.4	1.6	1.5	1.0
	PyT. Overhead	3.5	4.9	2.8	3.2	4.0	3.2
	Fut. Overhead	2.0	1.8	1.9	2.7	2.8	2.8
MI100	PyT. Jacob. (ms)	20.9	51.5	42.5	20.7	38.5	193.1
	Fut. Speedup	3.3	4.0	2.1	2.9	2.5	1.7
	PyT. Overhead	5.9	5.3	2.4	2.6	3.1	2.8
	Fut. Overhead	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 J	acob.	Fut	hark	nn.	LSTM	JAX	JAX(v	map)
A100	<b>D</b> <sub>0</sub>	45.4 ms			3.0		11.6	4.5		0.3
	$\mathbf{D}_1$	740	.1 ms		3.3		22.1	6.4		0.9
00	<b>D</b> <sub>0</sub>	89	.8 ms		2.6		4.0	_		_
Μ1	$\mathbf{D}_1$	1446	.9 ms		1.8		5.4	_		_
		Overheads								
		PyTorch	Futha	ırk	nn.LS	ТΜ	JAX	JAX(v	map)	
A100	<b>D</b> <sub>0</sub>	4.1	2	2.1		2.7	3.5		1.4	
	$\mathbf{D}_1$	4.3	3	5.9		2.2	3.7		0.8	
MI100	<b>D</b> <sub>0</sub>	5.0	4	.2		7.2	_		_	
	$\mathbf{D}_1$	7.9	3	5.9		6.6	_		_	
cuD	NN-b	ased refer	s to the	e (ma	anual)	tor	ch.nr	ı.LST	M libra	rv.

# **GPU Benchmarks - Depth and Memory Consumption**



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

- 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
```

Reverse-AD: Core Re-Write Rule

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

Extra Slides Motivating Example for AD Differentiating Loops & Time-Space Tradeoff Differentiating Map Loops with In-Place Updates Reverse-AD: Core Re-Write Rule

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

Extra Slides Motivating Example for AD Differentiating Loops & Time-Space Tradeoff Differentiating Map

Loops with In-Place Updates







C	luster	sizes
Ţ	1	=1
1	1+1+1+1	=4
1	1+1+1+1+	11 = 5











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

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

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

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

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

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

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

Reverse-AD: Core Re-Write Rule

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

Extra Slides Motivating Example for AD Differentiating Loops & Time-Space Tradeoff Differentiating Map Loops with In-Place Updates



- 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.

y = 2 **for** *i* = 0...*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.

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

$$2 \text{ let } ys_0 = \text{scratch}(n, \\ 3 \qquad \text{sizeOf}(y_0)) \\ 4 \text{ let } (y'', ys) = \\ 5 \text{ loop } (y, ys) = (y_0, ys_0) \\ 6 \text{ for } i = 0 \dots n - 1 \text{ do} \\ 7 \text{ let } ys[i] = y \\ 8 \text{ stm}_{Stoop} \\ 9 \text{ in } (y', ys) \\ 12 \text{ let } (y''', \overline{fvs_l}) = \\ 13 \text{ loop } (\overline{y}, \overline{fvs_l}) = (\overline{y''}, \overline{fvs_{l_0}}) \\ 14 \text{ for } i = n - 1 \dots 0 \text{ do} \\ 15 \text{ let } y = ys[i] \\ 16 \text{ stm}_{Stoop} \\ 17 \text{ stm}_{Stoop} \\ 18 \text{ in } (\overline{y'}, \overline{fvs'_l}) \\ 19 \text{ let } \overline{y_0} + \overline{y'''} \\ \end{array} \right)$$

١

- 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  $ys_0 =$ scratch(n, $sizeOf(y_0)$ ) 3 4 let (y'', ys) =5 **loop**  $(y, ys) = (y_0, ys_0)$ Primal Trace 6 **for** *i* = 0...*n* - 1 **do** let  $y_s[i] = y$ 8 stmsloop 9 **in** (y', ys)12 let  $(\overline{y'''}, \overline{fvs_l}) =$ 13 loop  $(\overline{y}, \overline{fvs_l}) = (\overline{y''}, \overline{fvs_{lo}})$ **14** for  $i = n - 1 \dots 0$  do 15 let y = ys[i]16 stmsloop Reverse Sweep 17 stms<sub>loop</sub> 18 in(y', fvs')19 **let**  $\overline{y_0} += y'''$ 

- 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  $ys_0 =$ scratch(n,3  $sizeOf(y_0)$ ) 4 let (y'', ys) =5 **loop**  $(y, ys) = (y_0, ys_0)$ 6 **for** *i* = 0...*n* - 1 **do** let  $y_s[i] = y$ 8 stms<sub>loop</sub> 9 **in** (y', ys)12 let  $(\overline{y'''}, \overline{fvs_l}) =$ 13 loop  $(\overline{y}, \overline{fvs_l}) = (\overline{y''}, fvs_{l_0})$ 14 **for**  $i = n - 1 \dots 0$  **do** 15 **let** y = ys[i] 16 stmsloop 17 stms<sub>loop</sub> 18  $in(y', fvs'_i)$ 19 let  $\overline{v_0}$  +=  $\overline{v'''}$ 

Primal Trace

- 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 
$$ys_0 = scratch(n, 
3 sizeOf(y_0))$$
  
4 let  $(y'', y_s) = 
5 loop  $(y, y_s) = (y_0, y_{s_0})$   
6 for  $i = 0 \dots n - 1$  do  
7 let  $ys[i] = y$   
8  $stms_{loop}$   
9 in  $(y', y_s)$   
12 let  $(\overline{y'''}, \overline{fv_{s_l}}) =$   
13 loop  $(\overline{y}, \overline{fv_{s_l}}) = (\overline{y''}, \overline{fv_{s_{l_0}}})$   
14 for  $i = n - 1 \dots 0$  do  
15 let  $y = ys[i]$   
16  $\underline{stms_{loop}}$   
17  $\underline{stms_{loop}}$   
18 in  $(\overline{y'}, \overline{fv_{s_l'}})$   
19 let  $\overline{y_0} + = \overline{y'''}$$ 

1 1

1

Primal Trace

١

- 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 
$$ys_0 = scratch(n, 
3 sizeOf(y_0))$$
  
4 let  $(y'', y_s) = 
5 loop  $(y, y_s) = (y_0, y_{s_0})$   
6 for  $i = 0 \dots n - 1$  do  
7 let  $ys[i] = y$   
8  $stms_{loop}$   
9 in  $(y', y_s)$   
12 let  $(y''', \bar{fvs_l}) =$   
13 loop  $(\bar{y}, \bar{fvs_l}) = (\bar{y''}, \bar{fvs_{l_0}})$   
14 for  $i = n - 1 \dots 0$  do  
15 let  $y = ys[i]$   
16  $\underline{stms_{loop}}$   
17  $\underline{stms_{loop}}$   
18 in  $(\bar{y'}, \bar{fvs'_l})$   
19 let  $\bar{y_0} + = \bar{y'''}$$ 

1

1

1

Primal Trace

١

- 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 
$$ys_0 = scratch(n, 
3 sizeOf(y_0))$$
  
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  $(\overline{y'''}, \overline{fys_l}) = (\overline{y''}, \overline{fys_{l_0}})$   
14 for  $i = n - 1 \dots 0$  do  
15 let  $y = ys[i]$   
16  $\overline{stms_{loop}}$   
17  $\overline{stms_{loop}}$   
18 in  $(\overline{y'}, \overline{fys'_l})$   
19 let  $\overline{y_0} + = \overline{y'''}$$ 

1

1

1 1

1

Primal Trace

١

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

- For the original loop, we save  $n^3$  versions of y on the tape.
- For the strip-mined loop, only 3*n* 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).

Reverse-AD: Core Re-Write Rule

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

Extra Slides Motivating Example for AD Differentiating Loops & Time-Space Tradeoff Differentiating Map Loops with In-Place Updates

### 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

\uparrow

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

\uparrow

forall i = 0 \dots n - 1
xs[i] = as[i] * bs[i]
```

Differentiating map is straightforward in the absence of free variables

let 
$$\overline{as}, \overline{bs} = map \ (\lambda a \ b \ \overline{x} \ \overline{a_0} \ \overline{b_0} \rightarrow$$
  
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  $\overline{xs} \ \overline{as_0} \ \overline{bs_0}$ 

Maps involving free variables are more complicated to differentiate

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

Naive approach: turn free variables into bound variables.

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

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 aray variable as[i] can be updated with an operation 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.

Reverse-AD: Core Re-Write Rule

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

#### Extra Slides

Motivating Example for AD Differentiating Loops & Time-Space Tradeoff Differentiating Map

Loops with In-Place Updates

### Handling of Loops with In-Place Updates

#### The brownian loop brindge 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</pre>
```

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 brindge 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</pre>
```

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 reverse: 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 brindge 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</pre>
```

#### 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 reverse: 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!