

# **Futhark:** Purely Functional GPU-programming with Nested Parallelism and in-place Array Updates

Troels Henriksen, Niels G. W. Serup, Martin Elsmann, Fritz Henglein, Cosmin Oancea

Presented by:  
-Zaid Qureshi

# Motivation

- GPUs are traditionally programmed using sequential programming languages
  - Requires expertise to exploit the parallelism provided by GPUs
- Functional programming languages provide parallelizable primitives (ie. map, reduce, scan)
  - But when compiled naively, their performance is very bad

# Futhark

- Purely Functional Array programming language for GPUs
  - To ease GPU programming
- Expresses computation/parallelism using basic and streaming second-order array combinators (SOACs)
- Type system that allows expression of race-free in-place updates
- Compiler implements partial flattening to allow for more parallelism without destroying memory access patterns

# Futhark Syntax

$e ::= k \mid v$	(Constant/Variable)
$(v_1, \dots, v_n)$	(Tuple)
$v_1 \odot v_2$	(Apply binary operator)
<b>if</b> $v_1$ <b>then</b> $e_2$ <b>else</b> $e_3$	
$v[v_1, \dots, v_n]$	(Array indexing)
$v \ v_1 \ \dots \ v_n$	(Function call)
<b>let</b> $(p_1, \dots, p_n) = e_1$ <b>in</b> $e_2$	(Let binding)
$v$ <b>with</b> $[v_1, \dots, v_n] \leftarrow v$	(In-place update)
<b>loop</b> $(pv)$ <b>for</b> $v < v$ <b>do</b> $e$	(Loop)
<b>iota</b> $v$	( $[0, \dots, v - 1]$ )
<b>replicate</b> $n \ v$	( $[v, \dots, v]$ of size $n$ )
<b>rearrange</b> $(\bar{k}) \ v$	(Rearrange dimensions)
<b>map</b> $l \ v_1 \ \dots \ v_n$	
<b>reduce</b> $l \ (v_1, \dots, v_n) \ v'_1 \ \dots \ v'_n$	
<b>scan</b> $l \ (v_1, \dots, v_n) \ v'_1 \ \dots \ v'_n$	
<b>stream</b> <b>_seq</b> $l \ (v_1, \dots, v_n) \ v'_1 \ \dots \ v'_m$	
<b>stream</b> <b>_seq</b> $l \ v'_1 \ \dots \ v'_m$	
<b>stream</b> <b>_red</b> $l_1 \ l_2 \ (v_1, \dots, v_n) \ v'_1 \ \dots \ v'_m$	

## Basic SOACs

$$\text{map } f [a_1, \dots, a_n] = [f a_1, \dots, f a_n]$$

$$\text{reduce } \oplus 0_{\oplus} [a_1, \dots, a_n] = 0_{\oplus} \oplus a_1 \oplus \dots \oplus a_n$$

$$\text{scan } \oplus 0_{\oplus} [a_1, \dots, a_n] = [a_1, \dots, a_1 \oplus \dots \oplus a_n]$$

## Example Futhark Code

```
fun main (matrix : [n][m]f32): ([n][m]f32, [n]f32) =  
  map ( $\lambda$ row : ([m]f32, f32)  $\rightarrow$   
    let row' = map ( $\lambda$ x : f32  $\rightarrow$  x+1.0) row  
    let s = reduce (+) 0 row  
    in (row',s))  
  matrix
```

# Example Futhark Code

INPUT: nxm matrix

```
fun main (matrix : [n][m]f32): ([n][m]f32, [n]f32) =  
  map ( $\lambda$ row : ([m]f32, f32)  $\rightarrow$   
    let row' = map ( $\lambda$ x : f32  $\rightarrow$  x+1.0) row  
    let s = reduce (+) 0 row  
    in (row',s))  
  matrix
```

# Example Futhark Code

INPUT:  $n \times m$  matrix

OUTPUT: tuple of  
 $n \times m$  matrix, array of size  $n$

```
fun main (matrix :  $[n][m]f32$ ): ( $[n][m]f32$ ,  $[n]f32$ ) =  
  map ( $\lambda$ row : ( $[m]f32$ ,  $f32$ )  $\rightarrow$   
    let row' = map ( $\lambda$ x :  $f32 \rightarrow x+1.0$ ) row  
    let s = reduce (+) 0 row  
    in (row',s))  
  matrix
```




# Example Futhark Code

```
fun main (matrix : ([n][m]f32)) : ([n][m]f32, [n]f32) =  
  map (λrow : ([m]f32, f32) →  
    let row' = map (λx : f32 → x+1.0) row  
    let s = reduce (+) 0 row  
    in (row',s))  
  matrix
```

Map over the rows  
of the matrix

# Example Futhark Code

```
fun main (matrix [n][m]f32): ([n][m]f32, [n]f32) =  
  map ( $\lambda$ row  ?)  $\rightarrow$   
    let row' = map ( $\lambda$ x : f32  $\rightarrow$  x+1.0) row  
    let s = reduce (+) 0 row  
    in (row',s)  
matrix
```

# Example Futhark Code

```
fun main (matrix : [n][m]f32): ([n][m]f32, [n]f32) =  
  map ( $\lambda$ row (f32 [m]f32)  $\rightarrow$   
    let map ( $\lambda$ x : f32  $\rightarrow$  x+1.0) row  
    let s = reduce (+) 0 row  
    in (row',s))  
  matrix
```

Get sum of row

# Example Futhark Code

```
fun main (matrix : [n][m]f32): ([n][m]f32, [n]f32) =  
  map ( $\lambda$ row : ([m]f32, f32)  $\rightarrow$   
    let (row', s) in (map ( $\lambda$ x : f32  $\rightarrow$  x+1.0) row  
    let s in (+) 0 row  
    in (row', s))  
  matrix
```

Return tuple of  
new row and sum

## Example Futhark Code

```
fun main (matrix : [n][m]f32): ([n][m]f32, [n]f32) =  
  map ( $\lambda$ row : ([m]f32, f32)  $\rightarrow$   
    let row' = map ( $\lambda$ x : f32  $\rightarrow$  x+1.0) row  
    let s = reduce (+) 0 row  
    in (row',s))  
  matrix
```

# Parallel operator *sFold* and Streaming Operators

$$\text{sFold } (\oplus) f (v_1 \# \dots \# v_k) = (f \ \epsilon) \oplus (f \ v_1) \oplus \dots \oplus (f \ v_k)$$

# - concat

$\epsilon$  - empty partition

# Parallel operator *sFold* and Streaming Operators

$$\text{sFold} (\oplus) f (v_1 \# \dots \# v_k) = (f \epsilon) \oplus (f v_1) \oplus \dots \oplus (f v_k)$$

$$\mathbf{stream\_map} f = \text{sFold} (\#) f$$

# - concat  
ε - empty partition

Applies *f* to each partition and then concatenates the resulting partitions

# Parallel operator *sFold* and Streaming Operators

$$\text{sFold } (\oplus) f (v_1 \# \dots \# v_k) = (f \epsilon) \oplus (f v_1) \oplus \dots \oplus (f v_k)$$

$$\mathbf{stream\_map} f = \text{sFold } (\#) f$$

# - concat  
ε - empty partition

Applies *f* to each partition and then concatenates the resulting partitions

$$\mathbf{stream\_red} (\oplus) f = \text{sFold } ((\oplus) * (\#)) f$$

Extends **stream\_map** by allowing each chunk to produce an additional output which is reduced in parallel



# Sequential Histogramming in Futhark

```
let counts = loop (counts = replicate k 0) for i < n do  
  let cluster = membership[i]  
  in counts with [cluster] ← counts[cluster] + 1
```

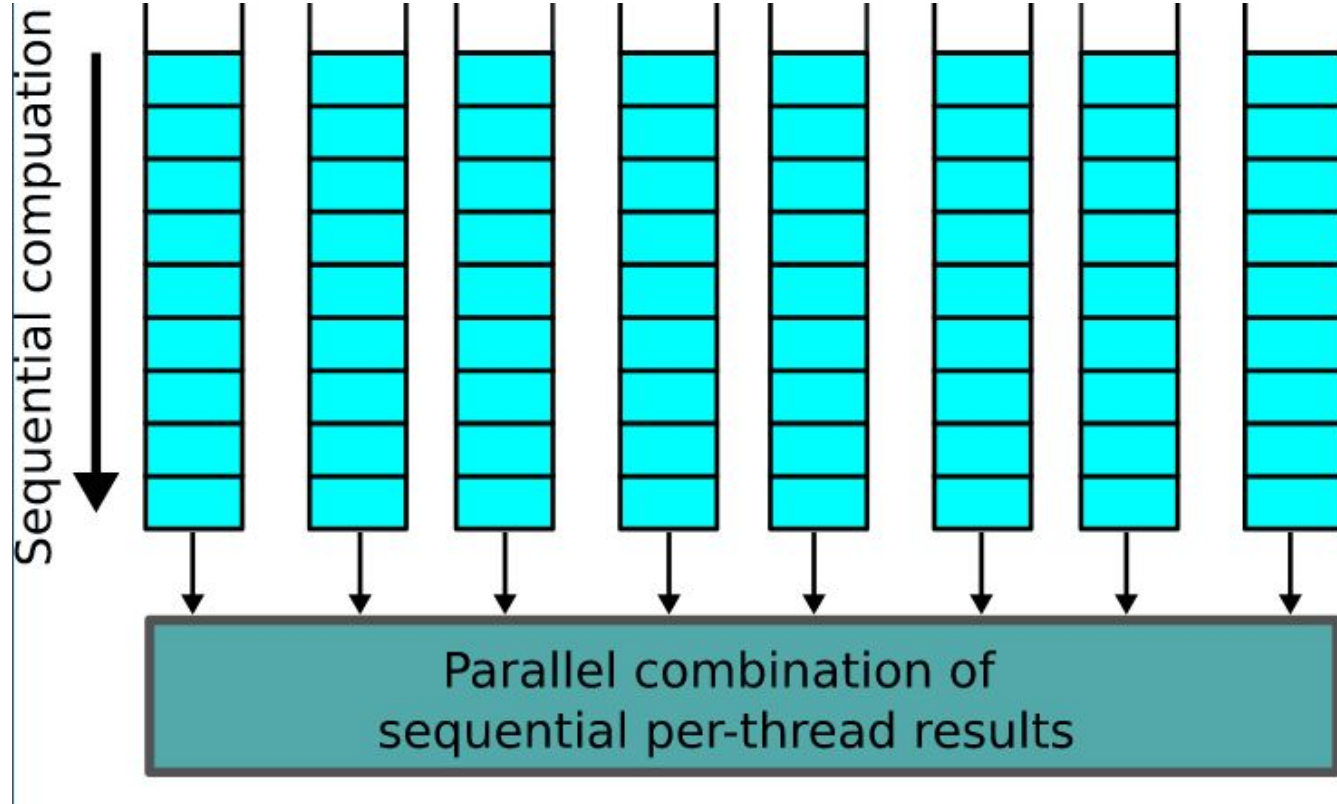
# Parallel Histogramming in Futhark

```
let increments =  
  map( $\lambda(\text{cluster}: \text{int}): [\text{k}]\text{int} \rightarrow$   
    let incr = replicate k 0 in let incr[cluster] = 1 in incr)  
  membership  
let counts =  
  reduce ( $\lambda(x: [\text{k}]\text{int}) (y: [\text{k}]\text{int}): [\text{k}]\text{int} \rightarrow$  map (+) x y)  
    (replicate k 0) increments
```

# Efficient Parallel Histogramming in Futhark

```
let counts = stream_red (map(+))
  (λ(acc: *[k]int) (chunk: [chunksize]int): [k]int →
    let res = loop (acc) for i < chunksize do
      let cluster = chunk[i]
      in acc with [cluster] ← acc[cluster]+1
    in res)
(replicate k 0) membership
```

# Efficient Parallel Histogramming in Futhark



# In-Place Updates and Uniqueness Types

```
fun modify (n: int) (a: *[n]int) (i: int) (x: [n]int): *[n]int =  
  a with [i] ← (a[i] + x[i])
```

- In purely functional languages array updates require copying array and updating copy (to avoid side effects)
- If it is known that the original array won't be used after the update, the update can occur in place

# In-Place Updates and Uniqueness Types

```
fun modify (n: int) (a: *[n]int) (i: int) (x: [n]int): *[n]int =  
  a with [i] ← (a[i] + x[i])
```

- Futhark has *uniqueness types* that allow programmer to specify function arguments that won't be referenced by the caller after the function call
  - The callee gains ownership of that argument
- An array is consumed when it is source of in-place update or when it is passed as a unique parameter.
- After the consumption point, the array or its aliases may not be used.
  - Type system checks this via aliasing rules

# Aliasing Rules

- Alias sets for values produced by SOACs are empty (new copies)
- Scalar read from an array does not alias its origin array (alias set not modified)
- Array slicing aliases origin array
- Function application:
  - If the result being returned is unique the alias set is empty
  - Otherwise the result aliases all non-unique parameters
- Other rules can be found in Figure 5 of the paper

# In-Place Update Checking

- Each expression  $e$  has a observed set of variables ( $O$ ) and a consumed set of variables ( $C$ )
  - the pair  $\langle C, O \rangle$  forms the occurrence trace for  $e$
- Inference rules used to check uniqueness and parameter consumption (Figure 6)

Sequence Judgement

$$\langle \mathcal{C}_1, \mathcal{O}_1 \rangle \gg \langle \mathcal{C}_2, \mathcal{O}_2 \rangle : \langle \mathcal{C}_3, \mathcal{O}_3 \rangle$$

Inference Rule

$$\frac{(\mathcal{O}_2 \cup \mathcal{C}_2) \cap \mathcal{C}_1 = \emptyset}{\langle \mathcal{C}_1, \mathcal{O}_1 \rangle \gg \langle \mathcal{C}_2, \mathcal{O}_2 \rangle : \langle \mathcal{C}_1 \cup \mathcal{C}_2, \mathcal{O}_1 \cup \mathcal{O}_2 \rangle} \quad (\text{OCCURENCE-SEQ})$$

If-then-else uniqueness inference rule

$$\frac{\begin{array}{l} v_1 \triangleright \langle \mathcal{C}_1, \mathcal{O}_1 \rangle \quad e_2 \triangleright \langle \mathcal{C}_2, \mathcal{O}_2 \rangle \quad e_3 \triangleright \langle \mathcal{C}_3, \mathcal{O}_3 \rangle \\ \langle \mathcal{C}_1, \mathcal{O}_1 \rangle \gg \langle \mathcal{C}_2, \mathcal{O}_2 \rangle : \langle \mathcal{C}'_2, \mathcal{O}'_2 \rangle \\ \langle \mathcal{C}_1, \mathcal{O}_1 \rangle \gg \langle \mathcal{C}_3, \mathcal{O}_3 \rangle : \langle \mathcal{C}'_3, \mathcal{O}'_3 \rangle \end{array}}{\text{if } v_1 \text{ then } e_2 \text{ else } e_3 \triangleright \langle \mathcal{C}'_2 \cup \mathcal{C}'_3, \mathcal{O}'_2 \cup \mathcal{O}'_3 \rangle} \quad (\text{SAFE-IF})$$



# In-Place Update Checking (Example)

```
let bs = map ( $\lambda(a) \rightarrow a$  with [0]  $\leftarrow$  2) as  
let d = iota m
```

This program passes as the function of the map consumes its parameter as

```
let cs = map ( $\lambda(i) \rightarrow d$  with [i]  $\leftarrow$  2) (iota n)
```

This program doesn't pass as it implies  $d$ , bound outside the function of the map, is consumed for every iteration of the map

# Streaming SOAC Fusion

$\mathcal{F}_1 : \text{map } f \bar{b} \implies \text{stream\_map } (\lambda(\bar{b}_c) \rightarrow \text{map } f \bar{b}_c) \bar{b}$   
 $\mathcal{F}_2 : \text{map } f \bar{b} \implies \text{stream\_seq } (\lambda(a, \bar{b}_c) \rightarrow (0, \text{map } f \bar{b}_c)) (0) \bar{b}$   
 $\mathcal{F}_3 : \text{reduce } \oplus \bar{e} \bar{b} \implies$   
 $\quad \text{stream\_red } \oplus (\lambda(\bar{a}, \bar{b}_c) \rightarrow \bar{a} \oplus \text{reduce } \oplus \bar{e} \bar{b}_c) (\bar{e}) \bar{b}$   
 $\mathcal{F}_4 : \text{reduce } \oplus \bar{e} \bar{b} \implies$   
 $\quad \text{stream\_seq } (\lambda(\bar{a}, \bar{b}_c) \rightarrow (\bar{a}) \oplus \text{reduce } \oplus \bar{e} \bar{b}_c) (\bar{e}) \bar{b}$   
 $\mathcal{F}_5 : \text{scan } \oplus \bar{e} \bar{b} \implies \text{stream\_seq}(\lambda(\bar{a}, \bar{b}_c) \rightarrow$   
 $\quad \text{let } \bar{x}_c = \text{scan } \oplus \bar{e} \bar{b}_c$   
 $\quad \text{let } \bar{y}_c = \text{map } (\bar{a} \oplus) \bar{x}_c$   
 $\quad \text{in } (\text{last}(\bar{y}_c), \bar{y}_c)) (\bar{e}) \bar{b}$

$\mathcal{F}_6 :$

$\text{let } (\bar{r}_1, \bar{x}, \bar{y}) =$   
 $\quad \text{stream\_red } \oplus f (\bar{e}_1) \bar{a}$   
 $\text{let } (\bar{r}_2, \bar{z}) =$   
 $\quad \text{stream\_red } \odot g (\bar{e}_2) \bar{x} \bar{b}$

$\mathcal{F}_7 :$

$\text{let } (\bar{r}_1, \bar{x}, \bar{y}) =$   
 $\quad \text{stream\_seq } f (\bar{e}_1) \bar{b} \implies$   
 $\text{let } (\bar{r}_2, \bar{z}) =$   
 $\quad \text{stream\_seq } g (\bar{e}_2) \bar{x} \bar{d}$

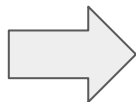
$\text{let } (\bar{r}_1, \bar{r}_2, \bar{x}, \bar{y}, \bar{z}) = \text{stream\_red}$   
 $\quad (\lambda(\bar{c}_1, \bar{d}_1, \bar{c}_2, \bar{d}_2) \rightarrow$   
 $\quad \quad (\bar{c}_1 \oplus \bar{c}_2, \bar{d}_1 \odot \bar{d}_2))$   
 $\quad (\lambda(\bar{e}_1, \bar{e}_2, \bar{a}_c, \bar{b}_c) \rightarrow$   
 $\quad \quad \text{let } (\bar{r}_1, \bar{x}_c, \bar{y}_c) = f \bar{e}_1 \bar{a}_c$   
 $\quad \quad \text{let } (\bar{r}_2, \bar{z}_c) = g \bar{e}_2 \bar{x}_c \bar{b}_c$   
 $\quad \quad \text{in } (\bar{r}_1, \bar{r}_2, \bar{x}_c, \bar{y}_c, \bar{z}_c))$   
 $\quad (\bar{e}_1, \bar{e}_2) \bar{a} \bar{b}$

$\text{let } (\bar{r}_1, \bar{r}_2, \bar{x}, \bar{y}, \bar{z}) = \text{stream\_seq}$   
 $\quad (\lambda(\bar{a}_1, \bar{a}_2, \bar{b}_c, \bar{d}_c) \rightarrow$   
 $\quad \quad \text{let } (\bar{r}_1, \bar{x}_c, \bar{y}_c) = f \bar{a}_1 \bar{b}_c$   
 $\quad \quad \text{let } (\bar{r}_2, \bar{z}_c) = g \bar{a}_2 \bar{x}_c \bar{d}_c$   
 $\quad \quad \text{in } (\bar{r}_1, \bar{r}_2, \bar{x}_c, \bar{y}_c, \bar{z}_c))$   
 $\quad (\bar{e}_1, \bar{e}_2) \bar{b} \bar{d}$

# Streaming SOAC Fusion Example

```
fun main(n: int): int =  
  let Y = stream_map ( $\lambda$ (iss: [m]int): [m]int  $\rightarrow$   
    let a = find (iss[0])  
    let t = map (g a) iss  
    let y = scan  $\odot$  0 t  
    in y)  
    (iota n)  
  let b = reduce (+) 0 Y  
  in b
```

(a) Program before fusion.



```
fun main(n: int): (int,int) =  
  stream_red (+) ( $\lambda$ (e1: int) (iss: [m]int): int  $\rightarrow$   
    let a = find(iss[0])  
    let t = map (g a) iss  
    let y = scan ( $\odot$ ) 0 t  
    let b = reduce (+) e1 y  
    in b)  
  (0) (iota n)
```

(b) Program after fusion at outer level.

# Moderate Flattening

- Flattening algorithm based on map-loop interchange and map distribution
- Attempt to exploit some top-level parallelism
  - Not seeking parallelism inside branches
  - Terminating map distribution when it would introduce irregular arrays

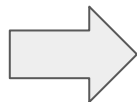
$$\text{map } f \circ \text{map } g \Rightarrow \text{map } (f \circ g)$$

# Moderate Flattening

- Flattening algorithm based on map-loop interchange and map distribution
- Attempt to exploit partial top-level parallelism
  - Not seeking parallelism inside branches
  - Terminating map distribution when it would introduce irregular arrays

$\text{map } f \circ \text{map } g \Rightarrow \text{map } (f \circ g)$

```
let bss: [m][m]i32 =  
  map (\(ps: [m]i32) (ps: [m]i32)  
->  
    loop (ws=ps) for i < n do  
      map (\w -> w * 2) ws)  
  pss
```

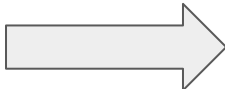


```
let bss: [m][m]i32 =  
  loop (wss=pss) for i < n do  
    map (\ws ->  
      map (\w -> w * 2) ws)  
  wss
```


# Locality of Reference Optimizations

- Naive translation of Flattened and Fused code can lead to bad memory access patterns
- Futhark compiler can optimize memory access patterns by transforming data

## Transpose:

`map (λxs → reduce (+) 0 xs) xss.`  `let xss' = as_column_major xss  
in map (λxs → reduce (+) 0 xs) xss'`

## Tiling:

`map (λp → stream_seq (λa (ps': [q]int) →  
g a ps') ps) ps`  `map (λp → stream_seq (λa (ps': [q]int) →  
let ps'' = local ps'  
in g a ps'') (0) ps) ps`

# Evaluation Methodology

- Tested with 2 GPUs
  - Nvidia GX 780
  - AMD W8100
- Generated OpenCL code is run on both GPUs
- Baseline implementations taken from benchmark suites

Benchmark	Dataset	
Backprop	Input layer size equal to $2^{20}$	Rodinia
CFD	<code>fvcorr.domn.193K</code>	
HotSpot	$1024 \times 1024$ ; 360 iterations	
K-means	<code>kdd_cup</code>	
LavaMD	<code>boxes1d=10</code>	
Myocyte	<code>workload=65536, xmax=3</code>	
NN	Default Rodinia dataset duplicated 20 times	
Pathfinder	Array of size $10^5$	
SRAD	$502 \times 458$ ; 100 iterations	
LocVolCalib	large dataset	
OptionPricing	large dataset	Parboil
MRI-Q	large dataset	
Crystal	Size 2000, degree 50	Accelerate
Fluid	$3000 \times 3000$ ; 20 iterations	
Mandelbrot	$4000 \times 4000$ ; 255 limit	
N-body	$N = 10^5$	

**Table 2:** Benchmark dataset configurations.

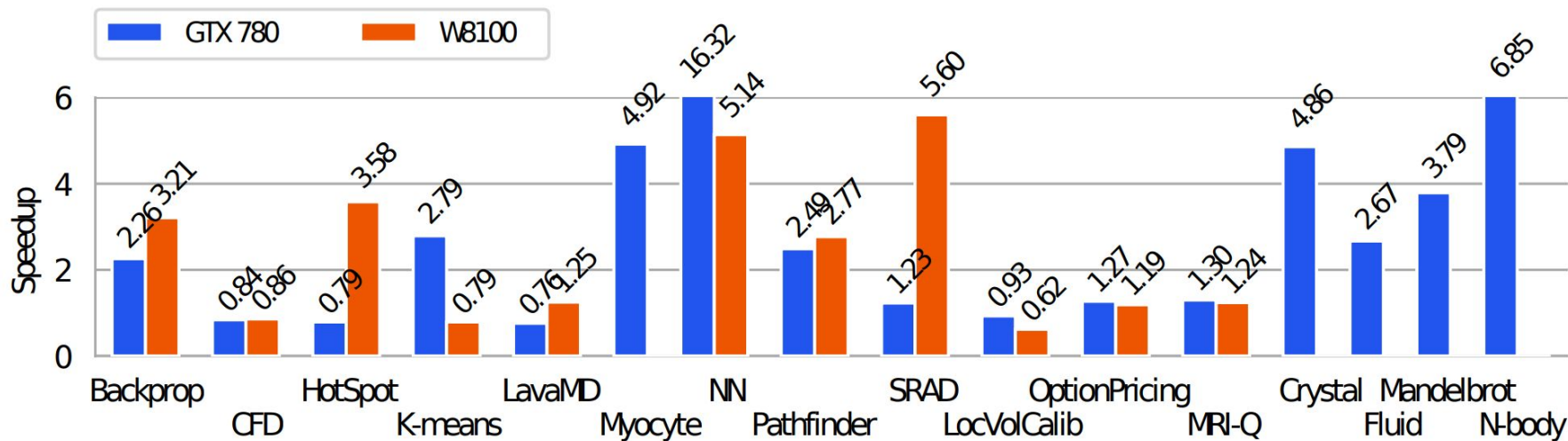
# Results

<b>Benchmark</b>	<b>NVIDIA GTX780</b>		<b>AMD W8100</b>	
	<b>Ref.</b>	<b>Futhark</b>	<b>Ref.</b>	<b>Futhark</b>
Backprop	46.9	20.7	41.5	12.9
CFD	1878.2	2235.9	3610.0	4177.5
HotSpot	35.9	45.3	260.4	72.6
K-means	1597.7	572.2	1216.1	1534.9
LavaMD	5.1	6.7	9.0	7.1
Myocyte	2733.6	555.4	—	2979.8
NN	178.9	11.0	193.2	37.6
Pathfinder	18.4	7.4	18.2	6.5
SRAD	19.9	16.1	195.1	34.8
LocVolCalib	1211.1	1293.2	3117.0	5015.8
OptionPricing	136.0	106.8	429.5	360.8
MRI-Q	20.2	15.5	17.9	14.3
Crystal	41.0	8.4	—	8.4
Fluid	268.7	100.4	—	221.8
Mandelbrot	30.8	8.1	—	14.8
N-body	613.2	89.5	—	269.8

**Table 1:** Average benchmark runtimes in milliseconds.

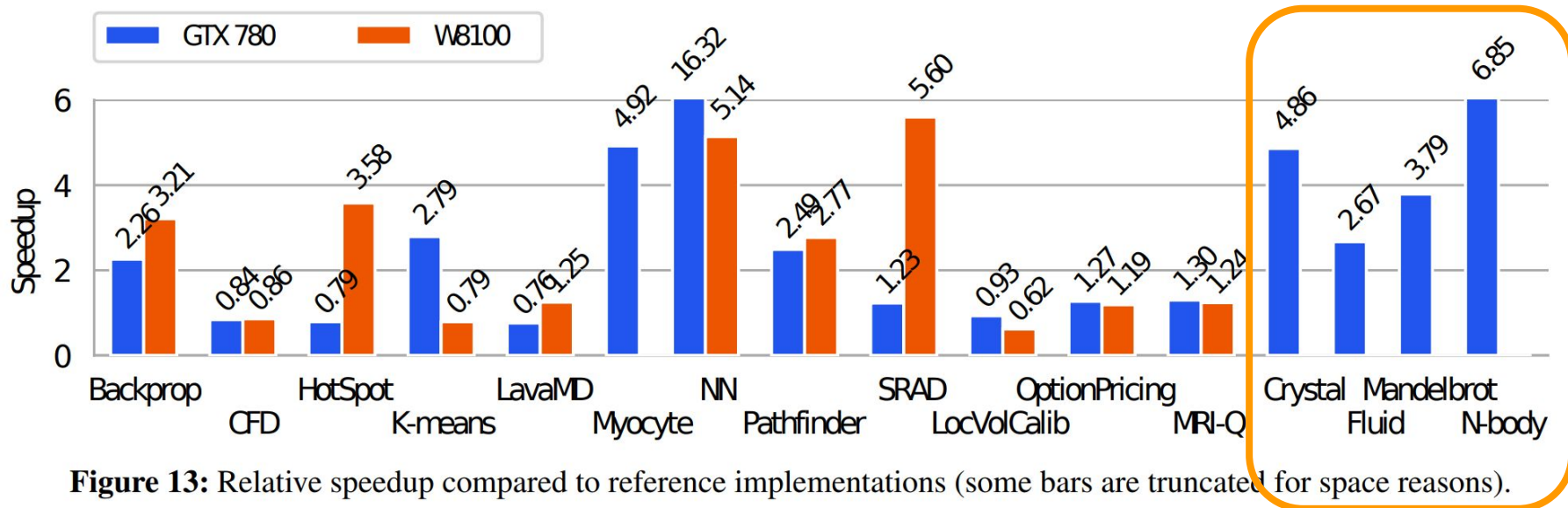


# Results



**Figure 13:** Relative speedup compared to reference implementations (some bars are truncated for space reasons).

# Results



**Futhark performs better than other functional programming environments for GPU due to higher level optimizations**

# Results

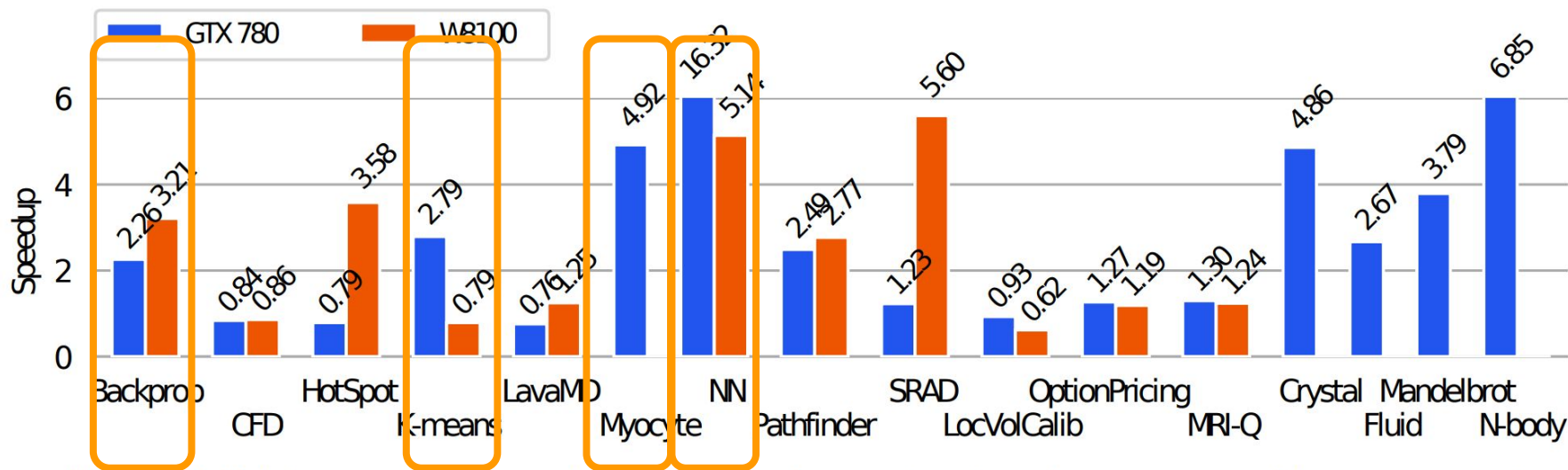


Figure 13: Relative speedup compared to reference implementations (some bars are truncated for space reasons).

**Rodinia doesn't implement all optimizations: sequential reductions (Backprop, NN), not parallelizing computation of new cluster centers (k-means), not coalescing all accesses (Myocyte)**

# Results

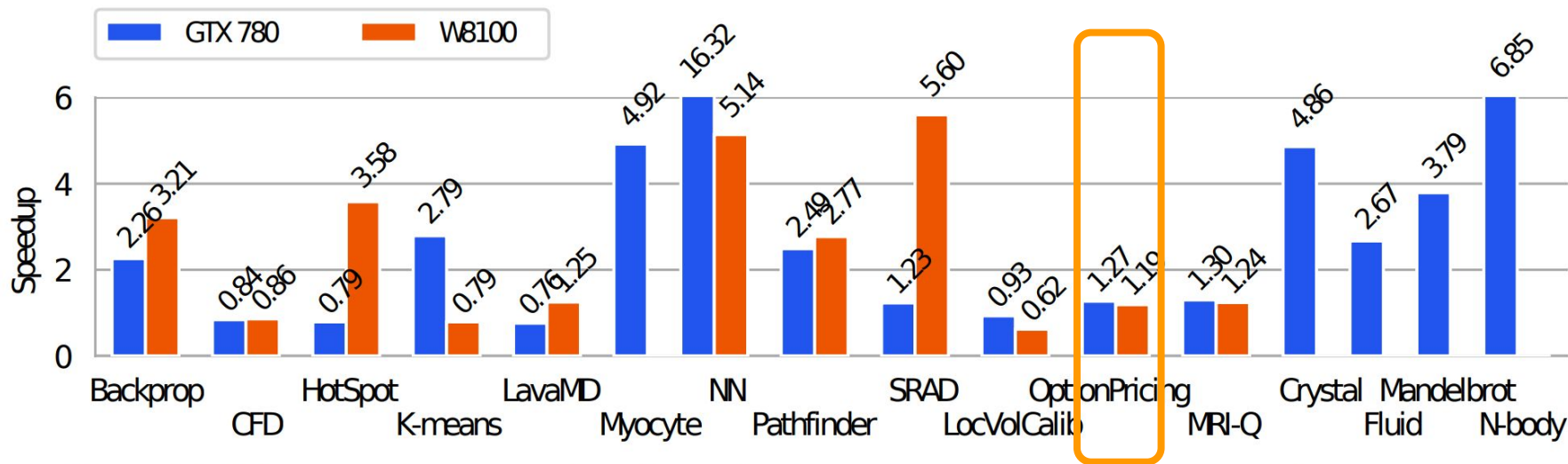


Figure 13: Relative speedup compared to reference implementations (some bars are truncated for space reasons).

**For OptionPricing, Futhark sequentializes excessive parallelism.**

# Results

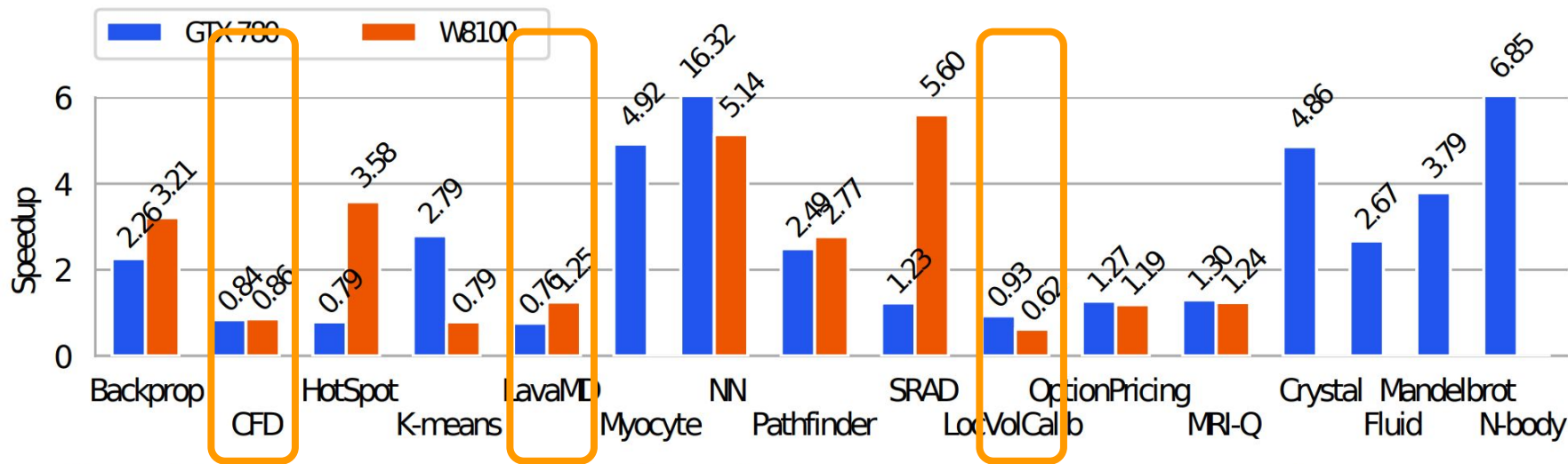


Figure 13: Relative speedup compared to reference implementations (some bars are truncated for space reasons).

**Furthark gets around 70-80% of the performance of hand-tuned code.**

# Impact of Optimizations

- SOAC Fusion
  - K-means (x1.42), LavaMD (x4.55), Myocyte (x1.66), SRAD (x1.21), Crystal (x10.1), LocVolCalid (x9.4)
  - Without fusion OptionPricing, N-body, and MRI-Q have too high memory requirements
- In-place Updates
  - K-means (x8.3), LocVolCalib (x1.7)
  - OptionPricing can't even be implemented without in-place updates
- Coalescing
  - K-means (x9.26), Myocyte (x4.2), OptionPricing (x8.79), LocVolCalib (x8.4)
- Tiling
  - LavaMD (x1.35), MRI-Q (x1.33), N-body (x2.29)

# Conclusion

## Pros:

- Futhark code is independent of the underlying hardware
- Futhark's type system allows expression of race-free in-place updates
- Optimizations done by compiler using higher level functions/reasoning
- Compiler implements partial flattening to allow for more parallelism without destroying memory access patterns
- Compiler can aggressively fuse and decompose code to best use available parallelism

## Cons:

- Requires rewrite of applications
- Although it does optimizations like flattening and fusion, Futhark's compiler can't optimize all the time
  - I.e. it can't convert inefficient histogramming to the efficient one
  - Still leaves a huge design space for the programmer to explore to write good performant code

Thank you!



Other slides

# Futhark Syntax

$e ::= k \mid v$	(Constant/Variable)
$(v_1, \dots, v_n)$	(Tuple)
$v_1 \odot v_2$	(Apply binary operator)
<b>if</b> $v_1$ <b>then</b> $e_2$ <b>else</b> $e_3$	
$v[v_1, \dots, v_n]$	(Array indexing)
$v \ v_1 \ \dots \ v_n$	(Function call)
<b>let</b> $(p_1, \dots, p_n) = e_1$ <b>in</b> $e_2$	(Let binding)
$v$ <b>with</b> $[v_1, \dots, v_n] \leftarrow v$	(In-place update)
<b>loop</b> $(pv)$ <b>for</b> $v < v$ <b>do</b> $e$	(Loop)
<b>iota</b> $v$	( $[0, \dots, v - 1]$ )
<b>replicate</b> $n \ v$	( $[v, \dots, v]$ of size $n$ )
<b>rearrange</b> $(\bar{k}) \ v$	(Rearrange dimensions)
<b>map</b> $l \ v_1 \ \dots \ v_n$	
<b>reduce</b> $l \ (v_1, \dots, v_n) \ v'_1 \ \dots \ v'_n$	
<b>scan</b> $l \ (v_1, \dots, v_n) \ v'_1 \ \dots \ v'_n$	
<b>stream</b> <b>_seq</b> $l \ (v_1, \dots, v_n) \ v'_1 \ \dots \ v'_m$	
<b>stream</b> <b>_seq</b> $l \ v'_1 \ \dots \ v'_m$	
<b>stream</b> <b>_red</b> $l_1 \ l_2 \ (v_1, \dots, v_n) \ v'_1 \ \dots \ v'_m$	

# Results

<b>Benchmark</b>	<b>NVIDIA GTX780</b>		<b>AMD W8100</b>	
	<b>Ref.</b>	<b>Futhark</b>	<b>Ref.</b>	<b>Futhark</b>
Backprop	46.9	20.7	41.5	12.9
CFD	1878.2	2235.9	3610.0	4177.5
HotSpot	35.9	45.3	260.4	72.6
K-means	1597.7	572.2	1216.1	1534.9
LavaMD	5.1	6.7	9.0	7.1
Myocyte	2733.6	555.4	—	2979.8
NN	178.9	11.0	193.2	37.6
Pathfinder	18.4	7.4	18.2	6.5
SRAD	19.9	16.1	195.1	34.8
LocVolCalib	1211.1	1293.2	3117.0	5015.8
OptionPricing	136.0	106.8	429.5	360.8
MRI-Q	20.2	15.5	17.9	14.3
Crystal	41.0	8.4	—	8.4
Fluid	268.7	100.4	—	221.8
Mandelbrot	30.8	8.1	—	14.8
N-body	613.2	89.5	—	269.8

**Table 1:** Average benchmark runtimes in milliseconds.

## Basic SOACs

$$\text{map } f [a_1, \dots, a_n] = [f a_1, \dots, f a_n]$$

$$\text{reduce } \oplus 0_{\oplus} [a_1, \dots, a_n] = 0_{\oplus} \oplus a_1 \oplus \dots \oplus a_n$$

$$\text{scan } \oplus 0_{\oplus} [a_1, \dots, a_n] = [a_1, \dots, a_1 \oplus \dots \oplus a_n]$$

Can be implemented with Parallel Operator *fold*

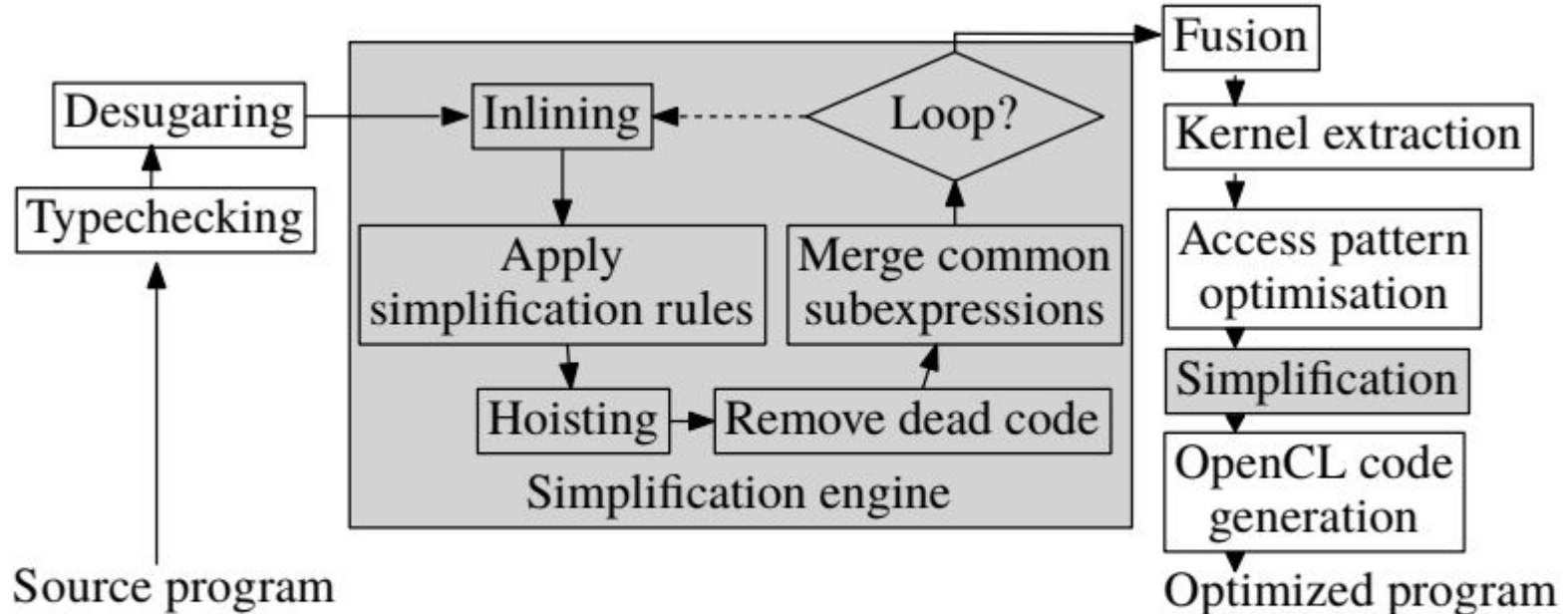
$$\text{fold } (\oplus, 0_{\oplus}) g [b_1, \dots, b_n] = 0_{\oplus} \oplus (g b_1) \oplus \dots \oplus (g b_n)$$

$$\begin{aligned}
\mathcal{F}1 : \text{map } f \bar{b} &\Longrightarrow \text{stream\_map } (\lambda(\bar{b}_c) \rightarrow \text{map } f \bar{b}_c) \bar{b} \\
\mathcal{F}2 : \text{map } f \bar{b} &\Longrightarrow \text{stream\_seq } (\lambda(a, \bar{b}_c) \rightarrow (0, \text{map } f \bar{b}_c)) (0) \bar{b} \\
\mathcal{F}3 : \text{reduce } \oplus \bar{e} \bar{b} &\Longrightarrow \\
&\quad \text{stream\_red } \oplus (\lambda(\bar{a}, \bar{b}_c) \rightarrow \bar{a} \oplus \text{reduce } \oplus \bar{e} \bar{b}_c) (\bar{e}) \bar{b} \\
\mathcal{F}4 : \text{reduce } \oplus \bar{e} \bar{b} &\Longrightarrow \\
&\quad \text{stream\_seq } (\lambda(\bar{a}, \bar{b}_c) \rightarrow (\bar{a}) \oplus \text{reduce } \oplus \bar{e} \bar{b}_c) (\bar{e}) \bar{b} \\
\mathcal{F}5 : \text{scan } \oplus \bar{e} \bar{b} &\Longrightarrow \text{stream\_seq} (\lambda(\bar{a}, \bar{b}_c) \rightarrow \\
&\quad \text{let } \bar{x}_c = \text{scan } \oplus \bar{e} \bar{b}_c \\
&\quad \text{let } \bar{y}_c = \text{map } (\bar{a} \oplus) \bar{x}_c \\
&\quad \text{in } (\text{last}(\bar{y}_c), \bar{y}_c)) (\bar{e}) \bar{b}
\end{aligned}$$

$$\begin{aligned}
\mathcal{F}6 : & \quad \text{let } (\bar{r}_1, \bar{r}_2, \bar{x}, \bar{y}, \bar{z}) = \text{stream\_red} \\
& \quad (\lambda(\bar{c}_1, \bar{d}_1, \bar{c}_2, \bar{d}_2) \rightarrow \\
& \quad (\bar{c}_1 \oplus \bar{c}_2, \bar{d}_1 \odot \bar{d}_2)) \\
& \quad \text{let } (\bar{r}_1, \bar{x}, \bar{y}) = \text{stream\_red } \oplus f (\bar{e}_1) \bar{a} \implies (\lambda(\bar{e}_1, \bar{e}_2, \bar{a}_c, \bar{b}_c) \rightarrow \\
& \quad \text{let } (\bar{r}_1, \bar{x}_c, \bar{y}_c) = f \bar{e}_1 \bar{a}_c \\
& \quad \text{let } (\bar{r}_2, \bar{z}_c) = g \bar{e}_2 \bar{x}_c \bar{b}_c \\
& \quad \text{in } (\bar{r}_1, \bar{r}_2, \bar{x}_c, \bar{y}_c, \bar{z}_c)) \\
& \quad \text{let } (\bar{r}_2, \bar{z}) = \text{stream\_red } \odot g (\bar{e}_2) \bar{x} \bar{b} \quad (\bar{e}_1, \bar{e}_2) \bar{a} \bar{b}
\end{aligned}$$

$$\begin{aligned}
\mathcal{F}7 : & \quad \text{let } (\bar{r}_1, \bar{r}_2, \bar{x}, \bar{y}, \bar{z}) = \text{stream\_seq} \\
& \quad (\lambda(\bar{a}_1, \bar{a}_2, \bar{b}_c, \bar{d}_c) \rightarrow \\
& \quad \text{let } (\bar{r}_1, \bar{x}_c, \bar{y}_c) = f \bar{a}_1 \bar{b}_c \\
& \quad \text{let } (\bar{r}_2, \bar{z}_c) = g \bar{a}_2 \bar{x}_c \bar{d}_c \\
& \quad \text{in } (\bar{r}_1, \bar{r}_2, \bar{x}_c, \bar{y}_c, \bar{z}_c)) \\
& \quad \text{let } (\bar{r}_1, \bar{x}, \bar{y}) = \text{stream\_seq } f (\bar{e}_1) \bar{b} \implies (\bar{e}_1, \bar{e}_2) \bar{b} \bar{d} \\
& \quad \text{let } (\bar{r}_2, \bar{z}) = \text{stream\_seq } g (\bar{e}_2) \bar{x} \bar{d}
\end{aligned}$$

# Futhark Compiler Architecture



**Figure 3:** Compiler architecture.