

Streaming in Accelerate

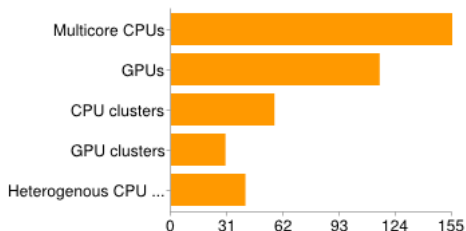
Frederik Meisner Madsen

DIKU
fmma@diku.dk

August 28, 2014

A short motivation

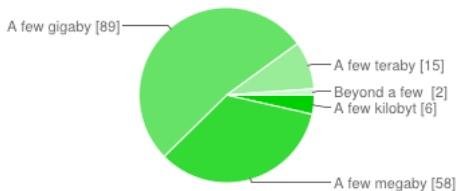
A recent survey on the future of array-oriented computing in Haskell by Manuel Chakravarty¹ - Interesting highlights: **Which type of high-performance hardware are you interested in using?**



¹<http://justtesting.org/post/70881852870/>

A short motivation

What is the typical size of your data sets?



A short motivation

- ▶ Conflict:
 - ▶ GPUs excel at computations on lots of data.
 - ▶ Not good for small problems due to low occupancy and overhead of GPU initialization.
 - ▶ Memory capacity is limited.
 - ▶ GeForce GTX 770: 4 GB device memory.
- ▶ Problem: Run out of memory fast.
- ▶ Solution: Manifest only when absolutely necessary.

Streaming

- ▶ Language design choice:
 - ▶ Manual streaming:
 - ▶ Easy for language implementer (do nothing).
 - ▶ Nightmare for programmer (manage buffers, scheduling, tied to specific platform).
 - ▶ Language-integrated streaming:
 - ▶ Nightmare for language implementer (to get right).
 - ▶ Easy for programmer.

NESL

- ▶ Forefather to Data Parallel Haskell.
- ▶ Based on (SIMD) vector-model, suitable for GPUs.
- ▶ Small language with formal (time) cost model, suitable for research.
 - ▶ Work and step.
- ▶ Most innovative feature: Vectorization of nested data parallelism.
 - ▶ Theoretical ideal asymptotic complexity.
 - ▶ With the right tricks: Scattered segment descriptors.
 - ▶ Space in order of exposed parallelism.
 - ▶ Naive matrix mult. requires $O(N^3)$ space.
 - ▶ Even moderate sized data sets can run into space problems.
 - ▶ Programmer should not be punished for exposing too much parallelism.

NESL + streams

- ▶ Goal: Language-integrated streaming.
- ▶ Two container types: Vectors $[\tau]$ and sequences $\{\sigma\}$.
- ▶ Sequences:
 - ▶ Semantically identical to vectors.
 - ▶ Processed one element at a time from first to last element.
 - ▶ No rewinding (time).
 - ▶ Reuse same memory for each element (space).
 - ▶ Backend is free to choose a chunk size.
 - ▶ Data parallelism.
 - ▶ CPU: Chunk size = 1.
 - ▶ GPU: Chunk size = 100.000.
- ▶ Sequence of vectors $\{[\tau]\}$:
 - ▶ Many important applications: Lines of a file, frames of a video.
 - ▶ Chunk must be able to grow dynamically to contain at least one vector.

NESL syntax (simplified)

$op :: \tau \rightarrow \tau$

	$+$	$:: (\text{Int}, \text{Int}) \rightarrow \text{Int}$
		\vdots
$e ::=$	$constant$	
	x	
	$\mathbf{let} \ x = e_1 \ \mathbf{in} \ e_2$	$mkvec_{\tau}^k :: (\overbrace{\tau, \dots, \tau}^k) \rightarrow [\tau]$
	$(e_1, \dots, e_k) \mid e.k$	$length_{\tau} :: [\tau] \rightarrow \text{Int}$
	$op \ e$	$!_{\tau} :: ([\tau], \text{Int}) \rightarrow \tau$
	$[e_1 : x \ \mathbf{in} \ e_2]$	$concat_{\tau} :: [[\tau]] \rightarrow [\tau]$
	$(= \text{map} (\lambda x. e_1) \ e_2)$	$partition_{\tau} :: ([\tau], [\text{Int}]) \rightarrow [[\tau]]$
	$scan$	$:: [\text{Int}] \rightarrow [\text{Int}]$
	sum	$:: [\text{Int}] \rightarrow \text{Int}$
		\vdots

NESL + streams syntax

$sop :: \sigma \rightarrow \sigma$

$e ::= \dots$
| $sop\ e$
| $\{e_1 : x\ \mathbf{in}\ e_2\}$

$mkseq_{\sigma}^k :: \overbrace{(\sigma, \dots, \sigma)}^k \rightarrow \{\sigma\}$
 $sconcat_{\sigma} :: \{\{\sigma\}\} \rightarrow \{\sigma\}$
 $flagpart_{\sigma} :: (\{\sigma\}, \{\mathbf{Bool}\}) \rightarrow \{\{\sigma\}\}$
 $sscan :: \{\mathbf{Int}\} \rightarrow \{\mathbf{Int}\}$
 $ssum :: \{\mathbf{Int}\} \rightarrow \mathbf{Int}$
 $tab :: \{\tau\} \rightarrow [\tau]$
 $seq :: [\tau] \rightarrow \{\tau\}$
 \vdots

NESL + streams

- ▶ Sequence operations can simulate almost all vector operations, except random-access and constant-time length:

$$\text{concat}_{\text{Int}}(e)$$
$$\equiv$$
$$\text{tab}_{\text{Int}}(\text{sconcat}_{\text{Int}}(\{\text{seq}_{\text{Int}}(x) : x \text{ in } \text{seq}_{[\text{Int}]}(e)\}))$$

- ▶ Same time complexity in cost model.
- ▶ Similarly, sequence comprehension can simulate vector comprehension.
- ▶ After eliminating redundant syntax, the only vector-related syntactic constructions are:

$(!_{\tau})$, length_{τ} , tab_{τ} and seq_{τ} .

NESL + streams

- ▶ Variables in sequence comprehension body cannot have sequence type:
 - ▶ Otherwise, sequence is reused once per element.
 - ▶ Easily checked by type system.
 - ▶ Workarounds:
 - ▶ Explicit tabulation using tab_{τ} .
 - ▶ Explicit recomputation by inlining.
 - ▶ Instead of silently making the choice in the compiler, resulting in two significantly different time/space complexities, the programmer is forced to make the choice.
- ▶ Vectorization:
 - ▶ Note: No vector-comprehension.
 - ▶ No scattered segment descriptors for sequences.
 - ▶ Type system forbids the cases where it would be needed.

NESL + streams

- ▶ Memory cost model.
 - ▶ Work and step analogue.
- ▶ Programming experience is almost as NESL.
 - ▶ Forced to make decisions about what should be vectors and what should be sequences.
 - ▶ ... but without a particular backend in mind.
 - ▶ Easy to reason about space.
- ▶ Open issues:
 - ▶ Feasible in practice?
 - ▶ Large constants in time complexity.
 - ▶ Schedulability / rate analysis.
 - ▶ Type system does not reject

$$\{sum(xs)\} ++ xs.$$

- ▶ Note that

$$xs ++ \{sum(xs)\}$$

is perfectly fine.

- ▶ Cost preservation theorem.

Accelerate

- ▶ DSL embedded in Haskell.
- ▶ Based on multi-dimensional array operations.
- ▶ Flat data parallelism (for now).
 - ▶ Regular nesting, the rows of a matrix all have the same length.
- ▶ GPU backend.

Accelerate

Scalar a = Array Z a

Vector a = Array (Z :: Int) a

Matrix a = Array (Z :: Int :: Int) a

generate :: $sh \rightarrow (sh \rightarrow a) \rightarrow \text{Array } sh \ a$

map :: $(a \rightarrow b) \rightarrow \text{Array } sh \ a \rightarrow \text{Array } sh \ b$

zipWith :: $(a \rightarrow b \rightarrow c) \rightarrow \text{Array } sh \ a \rightarrow \text{Array } sh \ b \rightarrow \text{Array } sh \ c$

scanl :: $(a \rightarrow a \rightarrow a) \rightarrow a \rightarrow \text{Vector } a \rightarrow \text{Vector } a$

fold :: $(a \rightarrow a \rightarrow a) \rightarrow a \rightarrow \text{Array } (sh :: Int) \ a \rightarrow \text{Array } sh \ a$

- ▶ Omitted: Everything wrapped in Acc.
 - ▶ Operations construct AST terms, $run :: \text{Acc } a \rightarrow a$.
 - ▶ All sharing is lost initially, recovered using De Bruijn indices.

Accelerate + streams

- ▶ What I hope to gain:
 - ▶ Data parallel streaming feasible in practice?
 - ▶ Contribution to Accelerate.
- ▶ Main challenge:
 - ▶ NESL + streams: Experimental toy language, designed from scratch with streaming in mind.
 - ▶ Accelerate + streams: Add streaming to real-world language.
- ▶ A multi-dimensional array in Accelerate is almost the same as a flat vector in NESL.
 - ▶ Both are fully manifest.
 - ▶ Data has identical representations in memory.
- ▶ Accelerate is similar to NESL.
 - ▶ Operations take array extents as additional arguments.
 - ▶ Specialized array operations (e.g. stencil).
 - ▶ Does not have nested vectors or vector-comprehensions.

Accelerate + streams - Take 1

- ▶ Like NESL + streams, introduce sequence container surface type.
- ▶ Unlike NESL + streams, do not eliminate existing vector operations.
 - ▶ Accelerate is already optimized for high performance.
 - ▶ Existing benchmarks should not become worse.
 - ▶ Breaks other backends.
 - ▶ Reimplementing everything defeats the purpose of using Accelerate in the first place.
- ▶ Sequence = Ordinary Haskell list.
 - ▶ $\{\text{Int}\} \simeq [\text{Scalar Int}] = [\text{Array Z Int}]$.
 - ▶ $\{\tau\} \simeq [\text{Array } (sh \text{ } \tau) \tau]$.
 - ▶ Get streaming from Haskell's lazy evaluation strategy.

Accelerate + streams - Take 1

type A = Array

mapStream :: (A sh a → A sh' b)
 → [A sh a] → [A sh' b]

zipWithStream :: (A sh a → A sh' b → A sh'' c)
 → [A sh a] → [A sh' b] → [A sh'' c]

foldStream :: (A sh a → A sh a → A sh a)
 → A sh a
 → [A sh a]
 → A sh a

toStream :: A (sh :: Int) a → [A sh a]

fromStream :: [A sh a] → (Vector sh, Vector a)

Accelerate CUDA backend

- ▶ Execution:

```
executeOpenAcc  ::  ExecOpenAcc aenv arrs  
                  →  aenv  
                  →  CIO arrs
```

- ▶ ExecOpenAcc *aenv arrs*: Executable AST.
 - ▶ Sharing recovery.
 - ▶ Fusion.
 - ▶ CUDA Code generation.
 - ▶ *aenv*:
 - ▶ Type-level list.
 - ▶ Sharing context from previous let bindings.
 - ▶ *arrs*: Result type (e.g. Vector Int).

Accelerate CUDA backend

- ▶ CIO = ReaderT Context (StateT State IO)
 - ▶ Context: Device properties and execution context.
 - ▶ State:
 - ▶ Host/device memory associations.
 - ▶ Compiled kernel object code.

Accelerate + streams - Take 1

- ▶ $executeOpenAcc$ (MapStream f acc) $aenv$:: CIO [A sh a]:

$$executeOpenAcc$$
 (MapStream f acc) $aenv$ =
$$do$$
 as \leftarrow $executeOpenAcc$ acc $aenv$
$$mapM$$
 ($executeOpenAfun$ f $aenv$) as

- ▶ Problem: Both state and IO monad are strict, all elements will manifest.
 - ▶ Lazy IO?
- ▶ Coming up with a solution - 4 failed approaches.

Accelerate + streams - Swapping types

- ▶ Use `[CIO arrs]` instead of `CIO [arrs]`?
- ▶ Define new function:

```
streamOpenAcc  :: ExecOpenAcc aenv [arrs]  
                → aenv  
                → [CIO arrs]
```

```
streamOpenAcc (MapStream f acc) aenv =  
  let s = streamOpenAcc acc aenv  
  in  map (>>= executeOpenAcc f aenv) s
```

- ▶ For sub-expression of stream type in `executeOpenAcc`, call `streamOpenAcc` instead of recursion.
- ▶ Problem: Sharing streams:
 - ▶ Add `[CIO arrs]` to sharing context - Each use recomputes.
 - ▶ Run stream and add `[arrs]` to sharing context - Tabulation.

Accelerate + streams - Look ahead

- ▶ Run streams immediately:
 - ▶ When a stream producer (*toStream*) is encountered:
 - ▶ Traverse AST to find it's consumers.
 - ▶ Feed all elements immediately.
 - ▶ In *fromStream* and *foldStream* nodes, store result somewhere.
 - ▶ When a *fromStream* or *foldStream* is encountered, simply fetch the stored result.
- ▶ Problems:
 - ▶ A seemingly irrelevant AST node may introduce a new bindings required by a later consumer.
 - ▶ *zipWithStream* requires two producers feeding elements in lock-step.
 - ▶ Traverse AST to find all producers in the same "loop" before starting to feed elements.
 - ▶ Many traversals, static analysis.

Accelerate + streams - Pipes

- ▶ Pipes is a popular library designed as a safe replacement to lazy IO.
- ▶ Combines effects, streaming and compositionality.
- ▶ Use

`Pipe aenv arrs CIO ()`

instead of

`aenv → CIO [arrs]?`

- ▶ Define new function:

`pipeOpenAcc` :: `ExecOpenAcc aenv [arrs]`
→ `Pipe aenv arrs CIO ()`

- ▶ Bind `arrs` instead of `[arrs]` in `aenv`. Feed `aenv` multiple times.

Accelerate + streams - Pipes

- ▶ Mapping closed functions:

$$\text{pipeOpenAcc (MapStream } f \text{ acc)} = \\ \text{pipeOpenAcc acc } \rightsquigarrow \text{mapPipe (execClosedAfun } f \text{)}$$

- ▶ Mapping open functions / let bindings:

$$\text{pipeOpenAcc (Alet } bnd \text{ bdy)} = \\ (\text{pipeOpenAcc } bnd \times \text{idPipe}) \rightsquigarrow \text{mapPipe pipeOpenAcc bdy}$$

- ▶ (\times) does not exist for Pipe.
 - ▶ A pipe may consume and produce in any order.

Accelerate + streams - Pipes

- ▶ Workarounds:

- ▶ Produce (Either a b) instead of (a, b) .
 - ▶ Context becomes a sum type:

Pipe $(\sum aenv)$ $arrs$ CIO $()$

- ▶ Buffering.
 - ▶ Upstream communication: Reject, request.
- ▶ Define a less general version of Pipe where a value is always produced directly after a value is consumed.
 - ▶ Better.

Accelerate + streams - Pipes

- ▶ Problem number 2:

$bnd :: \text{ExecOpenAcc } aenv [a]$
 $bdy :: \text{ExecOpenAcc } aenv b \quad (b \neq [-])$
 $executeOpenAcc (\text{Alet } bnd \ bdy) =??$

- ▶ Bind the pipe \rightarrow recomputation.
- ▶ Tabulate pipe, bind the result \rightarrow tabulation.
- ▶ Make everything a pipe.

Accelerate + streams - Concurrency

- ▶ Scheduling shared streams seems to be the core of the problem.
- ▶ Task parallelism is natural for stream transformers.
- ▶ Use (MVar *arrs*) instead of [*arrs*]?
- ▶ Execute stream operation:
 - ▶ Fork stream transformer. Take, compute, put loop.
- ▶ Problems:
 - ▶ Barrier synchronization for multiple consumers.
 - ▶ Synchronize State.
 - ▶ Synchronize communication with GPU.
 - ▶ ... and not just for streams. A binary array operation must wait if one operand is the result of a stream reduction.

Accelerate + streams - Take 1 conclusion

- ▶ Sharing streams is fundamentally different than sharing arrays.
 - ▶ Imposes restrictions on evaluation order.
- ▶ Streams as first-class “changes” all other types.
 - ▶ E.g. Int could be an integer available now, or an integer available in the future.
 - ▶ Not a problem in NESL + stream, compiler transforms everything to streams.
- ▶ Stream operations promise too much:
 - ▶ Partiality in the language.
 - ▶ Silent stream tabulation.

Accelerate + streams - Take 2

- ▶ Separate stream sharing contexts from array sharing context.
- ▶ New AST: Loop *lenv arrs*
 - ▶ Put all stream operations in Loop.
 - ▶ *lenv*: Loop sharing context (type-level list).
 - ▶ *arrs*: Result(s) of running the loop.
 - ▶ Allow closed loops in Accelerate AST:

$$\text{loop} :: \text{Loop } () \text{ arrs} \rightarrow \text{Acc arrs}$$

- ▶ Sequence = De Bruijn index into *lenv*.
 - ▶ $\{\text{Int}\} \simeq \text{Idx } \text{lenv } \text{Int}$.

data $\text{Idx } \text{env } a$ *where*

$$\text{ZeroIdx} :: \text{Idx } (\text{env}, a) a$$
$$\text{SuccIdx} :: \text{Idx } \text{env } a \rightarrow \text{Idx } (\text{env}, b) a$$

Accelerate + streams - Take 2

```
emptyLoop      :: Loop lenv ()  
mapStream     :: (A sh a → A sh' b)  
                → Idx lenv (A sh a)  
                → Loop (lenv, A sh' b) arrs  
                → Loop lenv arrs  
zipWithStream :: (A sh a → A sh' b → A sh'' c)  
                → Idx lenv (A sh a)  
                → Idx lenv (A sh' b)  
                → Loop (lenv, A sh'' c) arrs  
                → Loop lenv arrs
```

Accelerate + streams - Take 2

foldStream :: (A sh a → A sh a → A sh a)
→ A sh a
→ Idx lenv (A sh a)
→ Loop lenv arrs
→ Loop lenv (arrs, A sh a)

toStream :: A (sh :: Int) a
→ Loop (lenv, A sh a) arrs
→ Loop lenv arrs

fromStream :: Idx sh a
→ Loop lenv arrs
→ Loop lenv (arrs, (Vector sh, Vector a))

Accelerate + streams - Take 2

- ▶ Loop terminates when the first producer is exhausted (like *zip* in Haskell or vector-comprehension in NESL).
- ▶ Status: Reference interpreter works, confident that CUDA backend will work too.
- ▶ Difference from take 1:
 - ▶ Same operations.
 - ▶ Cannot let bind new arrays in Loop.
 - ▶ Move array let-bindings to just before Loop:
 - ▶ If not possible, array is a reduction of one of the streams of the loop.
 - ▶ Move to after Loop.
 - ▶ If not possible, a stream of the loop depends on a reduction of the loop → exactly when partiality or silent stream tabulation was required in take 1.
- ▶ This is as far as I got.

Accelerate + streams - Take 2

$$\sum_{i=1}^n \log i$$

```
iota :: Int -> Acc (Vector Int)
iota n = generate (index1 (constant n)) unindex1

-- Take 1
logsum :: Int -> Acc (Scalar Float)
logsum n = foldStream (zipWith (+)) (use (fromList Z [0.0]))
    $ mapStream (map (log . fromIntegral . (+1)))
    $ toStream (iota n)

-- Take 2
logsum :: Int -> Acc (Scalar Float)
logsum n = asnd $ loop
    $ toStream (iota n)
    $ mapStream (map (log . fromIntegral . (+1))) ZeroIdx
    $ foldStream (zipWith (+)) (use (fromList Z [0.0])) ZeroIdx
    $ emptyLoop
```


Accelerate + streams - Take 2 conclusion

- ▶ Essentially a sub-language for sequences.
- ▶ Many open questions:
 - ▶ Surface language with sharing recovery.
 - ▶ Should be easy enough.
 - ▶ More stream operations.
 - ▶ Scan, filter.
 - ▶ Generalize to chunk size > 1 .
 - ▶ Treat scalars (fixed-buffer) different from non-scalars (dynamic buffer) in loop context.
 - ▶ Variable number of elements produced, loops become dataflow networks.
 - ▶ Nested data parallelism.
 - ▶ Lifted loops.
 - ▶ Nested loops.

Conclusion

- ▶ For GPUs, streaming is a necessity.
 - ▶ Language-integrated streaming is preferable in high-level languages.
- ▶ Language-integrated first-class streaming is all or nothing.
 - ▶ Nothing is a stream or everything is a stream.
 - ▶ NESL: Everything is lifted to sequence space.
 - ▶ Accelerate: Streams are confined to a sub-language.
 - ▶ Closed non-suspendable loop. No communication across loops.
- ▶ Data parallel streaming feasible in practice?
 - ▶ Perhaps, but difficult to get right.
 - ▶ No benchmarks yet.