

FCL: A low-level functional GPU language

(Work in progress)

Martin Dybdal
dybber@dybber.dk

DIKU
University of Copenhagen

7 July 2016

Joint work w. Mary Sheeran, Joel Svensson and Martin Elsman

Agenda

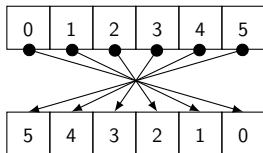
- ▶ Why a new low-level language for GPU computing?
- ▶ FCL by example
- ▶ Performance results
- ▶ Formalising FCL (just the highlights)
- ▶ Future work

Why a new low-level language for GPU computing?

- ▶ Composability and agility
- ▶ Built in fusion, user-controlled
- ▶ Compare programs/algorithms
- ▶ Predictability, no black box wrt. performance
- ▶ Ability to optimize
- ▶ A quest for a GPU language for algorithms researchers
- ▶ Intermediate language for optimizing compilers

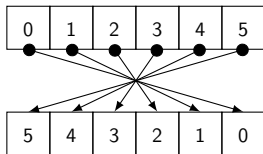
Reverse: A simple example

```
sig reverse : [a] -> [a]
fun reverse arr =
  let n = length arr
  in generate n (fn i => index arr (n - i - 1))
```



Reverse: A simple example

```
sig reverse : [a] -> [a]
fun reverse arr =
  let n = length arr
  in generate n (fn i => index arr (n - i - 1))
```



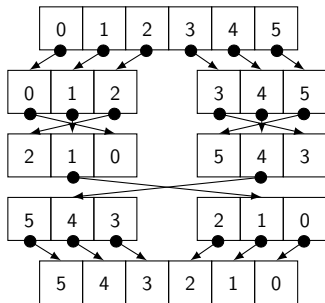
How is this executed on a GPU?

- ▶ Sequential in a single thread?
- ▶ Performed in a single block?
- ▶ Performed among threads in a grid?

Reverse: Distributed

```
sig revBlock : [a] -> [a]<block>  
fun revBlock arr = push <block> (reverse arr)
```

```
sig revDistribute : int -> [a] -> [a]<grid>  
fun revDistribute chunkSize arr =  
  splitUp chunkSize arr  
  |> map reverseBlock  
  |> reverse  
  |> concat chunkSize
```



Reverse: Generated OpenCL

```
sig revKernel : [int] -> [int]<grid>
kernel revKernel arr = revDistribute #BlockSize arr
  config #BlockSize = 256
```

```
__kernel void revKernel(__global int* arrInput_0, int lenInput_1,
                        __global int* arrOutput_3) {
    int n_2 = lenInput_1 / 256;
    int blocksQ_5 = n_2 / get_num_groups(0);
    for (int i_6 = 0; i_6 < blocksQ_5; i_6++) {
        int j_8 = (get_group_id(0) * blocksQ_5) + i_6;
        arrOutput_3[((j_8 * 256) + get_local_id(0))] =
            arrInput_0 [((((n_2 - j_8) - 1) * 256) + ((256 - get_local_id(0)) - 1))];
        barrier(CLK_LOCAL_MEM_FENCE);
    }
    if (get_group_id(0) < (n_2 % get_num_groups(0))) {
        int j_16 = (get_num_groups(0) * blocksQ_5) + get_group_id(0);
        arrOutput_3[((j_16 * 256) + get_local_id(0))] =
            arrInput_0 [((((n_2 - j_16) - 1) * 256) + ((256 - get_local_id(0)) - 1))];
        barrier(CLK_LOCAL_MEM_FENCE);
    }
}
```

Two array types

(from Obsidian)

- ▶ Pull arrays: for organizing computation, indexable no concatenation

`[int], [bool], [[int]]`

- ▶ Push arrays: for writing to memory, non indexable, supports concatenation

`[int]<thread>, [int]<block>, [int]<grid>`

- ▶ Both supporting fusion

Array types in reverse

```
sig revDistribute : int -> [a] -> [a]<grid>
fun revDistribute chunkSize arr =
  splitUp chunkSize arr      -- [[a]]
  |> map reverseBlock        -- [[a]<block>]
  |> reverse                  -- [[a]<block>]
  |> concat chunkSize        -- [a]<grid>

sig push : <lvl> -> [a] -> [a]<lvl>
sig splitUp : int -> [a] -> [[a]]
sig concat : int -> [[a]<lvl>] -> [a]<1+lvl>
```

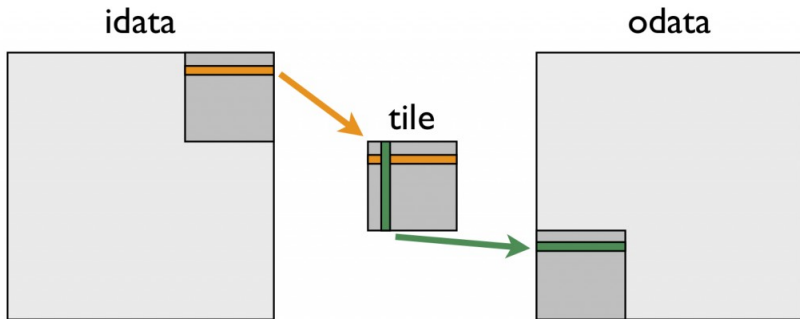
Matrix transposition

A basic approach

```
sig transpose : int -> int -> [a] -> [a]
fun transpose cols rows elems =
  generate (cols * rows)
    (fn n =>
      let i = n / rows in
      let j = n % rows
      in index elems (j * rows + i))
```

Matrix transposition

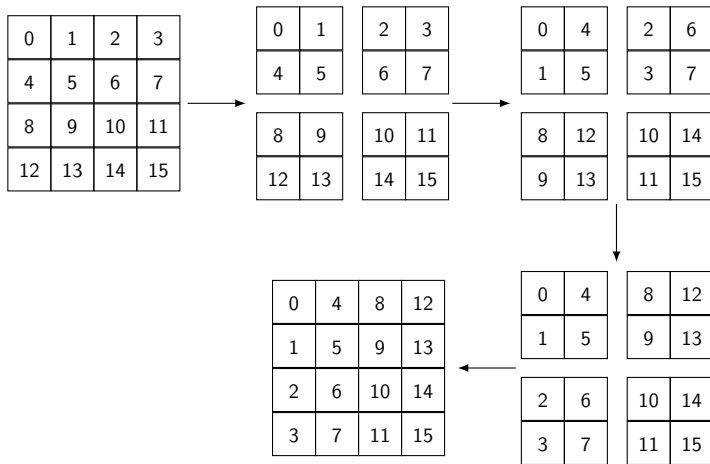
Tiled



(Figure by NVIDIA)

Matrix transposition

Tiled



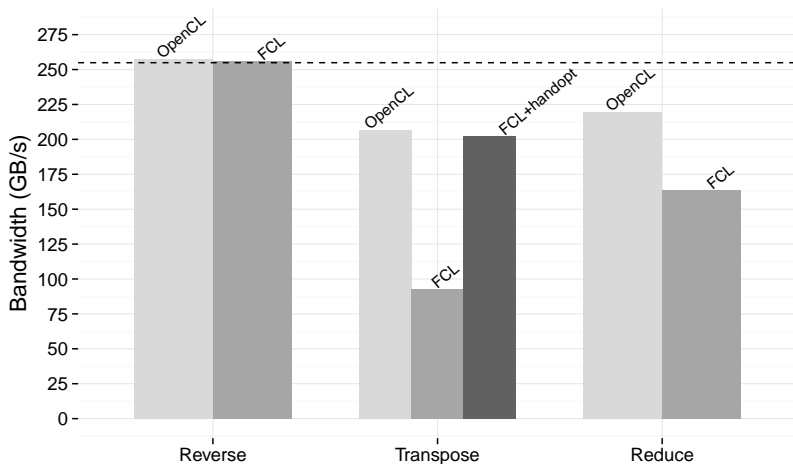
Matrix transposition

Tiled, using shared-memory

```
sig transposeTiled : int -> int -> int -> [a] -> [a]<grid>
fun transposeTiled tileDim cols rows mat =
  let n = cols / tileDim
      m = rows / tileDim
  in split2DGrid tileDim cols n m mat
      |> map (force . push <block>)
      |> map (transpose tileDim tileDim)
      |> transpose n m
      |> map (push <block>)
      |> concat2DGrid tileDim n rows

sig force : [a]<lvl> -> [a]
```

Performance



NVIDIA GeForce GTX 780 Ti (2880 cores, 875 Mhz, 3GB GDDR5)
Peak bandwidth 254.90 GB/s (measured)

Highlights from formalism

- ▶ Polymorphic type system, restricting e.g. the valid nesting of arrays
- ▶ Dynamic semantics, explaining mapping to threads, blocks, warps and grids
- ▶ Formalism informed our implementation
- ▶ Abstracts away from block/warp-virtualization
- ▶ Classic type safety properties

Built-in operators

$\text{lengthPull} : [\alpha] \rightarrow \text{int}$

$\text{lengthPush} : [\alpha]\langle l \rangle \rightarrow \text{int}$

$\text{mapPull} : (\alpha \rightarrow \beta) \rightarrow [\alpha] \rightarrow [\beta]$

$\text{mapPush} : (\alpha \rightarrow \beta) \rightarrow [\alpha]\langle l \rangle \rightarrow [\beta]\langle l \rangle$

$\text{generate} : \text{int} \rightarrow (\text{int} \rightarrow \alpha) \rightarrow [\alpha]$

$\text{index} : [\alpha] \rightarrow \text{int} \rightarrow \alpha$

$\text{push} : \langle l \rangle \rightarrow [\alpha] \rightarrow [\alpha]\langle l \rangle$

$\text{force} : [\alpha]\langle l \rangle \rightarrow [\alpha]$

$\text{concat} : \text{int} \rightarrow [[\alpha]\langle l \rangle] \rightarrow [\alpha]\langle 1 + l \rangle$

$\text{while} : ([\alpha] \rightarrow \text{bool}) \rightarrow ([\alpha] \rightarrow [\alpha]\langle l \rangle) \rightarrow [\alpha]\langle l \rangle \rightarrow [\alpha]$

Future work on FCL

- ▶ Multi-dimensional arrays
- ▶ Tracking communication costs in semantics
- ▶ Generalize hierarchy and mapping
(e.g. ability to add layers, like multiple GPUs)

Summary

- ▶ I argue that we need to focus on performance reasoning and cost-models
- ▶ Nested arrays, but only non-nested arrays can be materialized.
- ▶ Level hierarchy controls mapping to sequential/parallel loops
- ▶ FCL is very much work in progress

References



Obsidian: A domain specific embedded language for parallel programming of graphics processor

Joel Svensson, Mary Sheeran, Koen Claessen, 2011

IFL'11



Compiling a Subset of APL Into a Typed Intermediate Language.

Martin Elsman and Martin Dybdal, 2014

ARRAY'14

Thank you

Martin Dybdal

Ph.D. student

DIKU, University of Copenhagen

dybber@dybber.dk

FCL is available at: <http://github.com/dybber/fcl>

More future work on FCL

- ▶ Manual memory-management
- ▶ Sequential loops
- ▶ Larger examples
- ▶ Host-code generation
- ▶ Use as backend for our APL-compiler (TAIL)

Reduction

```
sig halve : [a] -> ([a], [a])
```

```
sig zipWith : (a -> b -> c) -> [a] -> [b] -> [c]
```

```
sig step : <lvl> -> (a -> a -> a) -> [a] -> [a]<lvl>
```

```
fun step <lvl> f arr =
```

```
  let x = halve arr
```

```
  in push <lvl> (zipWith f (fst x) (snd x))
```

0	1	2	3	4	5	6	7	8	9	10	11
---	---	---	---	---	---	---	---	---	---	----	----



0	1	2	3	4	5
6	7	8	9	10	11

+



6	8	10	12	14	16
---	---	----	----	----	----

Reduction

```
sig red : <lvl> -> (a -> a -> a) -> [a] -> [a]<lvl>
fun red <lvl> f arr =
  while (fn arr => 1 != lengthPull arr)
    (step <lvl> f)
    (step <lvl> f arr)
|> push <lvl>
```

Reduction

OpenCL

```
__kernel void reduceAdd(__local volatile uchar* sbase,
                        __global int* arrInput_0, int lenInput_1,
                        __global int* arrOutput_2) {
    int ub_3 = lenInput_1 >> 8;
    int blocksQ_4 = ub_3 / get_num_groups(0);
    for (int i_5 = 0; i_5 < blocksQ_4; i_5++) {
        int j_7 = (get_group_id(0) * blocksQ_4) + i_5;
        __local volatile int* arr_12 = (__local volatile int*) (sbase + 0);
        arr_12[get_local_id(0)] = arrInput_0 [((j_7 * 256) + get_local_id(0))]
                                + arrInput_0 [((j_7 * 256) + (get_local_id(0) + 128))];
        barrier(CLK_LOCAL_MEM_FENCE);
        if (get_local_id(0) < 64) {
            arr_12[get_local_id(0)] = arr_12 [get_local_id(0)] + arr_12 [(get_local_id(0) + 64)];
        }
        barrier(CLK_LOCAL_MEM_FENCE);
        if (get_local_id(0) < 32) {
            arr_12[get_local_id(0)] = arr_12 [get_local_id(0)] + arr_12 [(get_local_id(0) + 32)];
        }
        barrier(CLK_LOCAL_MEM_FENCE);
        if (get_local_id(0) < 16) {
            arr_12[get_local_id(0)] = arr_12 [get_local_id(0)] + arr_12 [(get_local_id(0) + 16)];
        }
        barrier(CLK_LOCAL_MEM_FENCE);
        if (get_local_id(0) < 8) {
            arr_12[get_local_id(0)] = arr_12 [get_local_id(0)] + arr_12 [(get_local_id(0) + 8)];
        }
        barrier(CLK_LOCAL_MEM_FENCE);
        if (get_local_id(0) < 4) {
            arr_12[get_local_id(0)] = arr_12 [get_local_id(0)] + arr_12 [(get_local_id(0) + 4)];
        }
        barrier(CLK_LOCAL_MEM_FENCE);
        if (get_local_id(0) < 2) {
            arr_12[get_local_id(0)] = arr_12 [get_local_id(0)] + arr_12 [(get_local_id(0) + 2)];
```