

DEPARTMENT OF COMPUTER SCIENCE AT THE UNIVERSITY OF COPENHAGEN

Bachelors project

Auto-tuning Futhark

Simon Rotendahl (mpx651) & Carl Mathias Graae Larsen (pwh334)
{*simon, cala*}@di.ku.dk

Contents

1	Introduction	1
2	Background	1
2.1	Futhark	1
2.2	Flattening	1
2.3	Incremental flattening	1
3	Design	1
A	Code examples	2
A.1	Matrix Multiplication - CUDA for GPUs [prog-guide-cuda]	2

May 2019

Abstract

Stuff

1 Introduction

Stuff

2 Background

2.1 Futhark

A common way to increase computer performance, is to increase the capacity for parallelism. For practical usage, however, this is difficult to implement, due to low-level GPU-specific languages requiring domain specific knowledge to make full use of that capacity. A vast amount of work has gone into transforming high-level hardware-agnostic code into these low-level GPU-specific languages [**inc-flat**].

The programming language **Futhark** aims to solve this problem. The creator of Futhark writes the purpose, of the language, nicely on the home page for the language *"Because it's nicer than writing CUDA or OpenCL by hand!"* [**futhark-home**]. On the same page, Futhark is described, more precisely, as *"a statically typed, data-parallel, and purely functional array language"*, but better than a description, is an example:

```
1 let dotprod [n] (xs: [n]f32) (ys: [n]f32): f32 =
2   reduce (+) 0f32 (map2 (*) xs ys)
3
4 let main [n][m][p] (xss: [n][m]f32) (yss: [m][p]
5   f32): [n][p]f32 =
6   map (\xs -> map (dotprod xs) (transpose yss))
7   xss
```

Figure 1: Matrix-matrix multiplication in Futhark [**ppopp**]

A Futhark program for matrix-matrix multiplication can be seen in figure 1, the syntax is similar to languages such as ML, and Haskell. It is a good example of how Futhark differs from CUDA or OpenCL (we would have liked to include an example of CUDA, but it was too long, so see A.1 for that). It allows the programmer to write efficient parallel code, without all the domain specific knowledge regarding massively parallel systems.

2.2 Flattening

It is difficult to exploit nested data-parallelism. An approach to this problem is flattening [**flat**]. The aim is to transform nested parallelism, into one-level parallelism. To see flattening, let's flatten a nested collection; let $A = [[1, 2], [3, 4, 5], [6]]$, this can be flattened into two collections (here represented as a tuple, but it could be a different data-structure) ($A_{\text{shape}} = [2, 3, 1]$, $A_{\text{val}} = [1, 2, 3, 4, 5, 6]$). This technique can be applied recursively to any depth, a depth of n would give n shape arrays, for example A would be represented as $A = ([3], (A_{\text{shape}}, A_{\text{val}})) \rightarrow ([3], ([2, 3, 1], [1, 2, 3, 4, 5, 6]))$

2.3 Incremental flattening

A principal critical for this code transformation from low-level GPU languages to Futhark, is flattening.

3 Design

A Code examples

A.1 Matrix Multiplication - CUDA for GPUs [prog-guide-cuda]

```
1 __global__ void Muld(float* A, float* B, int wA, int wB, float* C)
2 {
3     int bx = blockIdx.x;
4     int by = blockIdx.y;
5     int tx = threadIdx.x;
6     int ty = threadIdx.y;
7
8     int aBegin = wA * BLOCK_SIZE * by;
9     int aEnd   = aBegin + wA - 1;
10    int aStep   = BLOCK_SIZE;
11    int bBegin = BLOCK_SIZE * bx;
12    int bStep   = BLOCK_SIZE * wB;
13
14    float Csub = 0;
15    for (int a = aBegin, b = bBegin;
16         a <= aEnd;
17         a += aStep,
18         b += bStep) {
19
20        __shared__ float As[BLOCK_SIZE][BLOCK_SIZE];
21        __shared__ float Bs[BLOCK_SIZE][BLOCK_SIZE];
22
23        As[ty][tx] = A[a + wA * ty + tx];
24        Bs[ty][tx] = B[b + wB * ty + tx];
25
26        __syncthreads();
27
28        for (int k = 0; k < BLOCK_SIZE; ++k)
29            Csub += As[ty][k] * Bs[k][tx];
30    __syncthreads();
31    }
32
33    int c = wB * BLOCK_SIZE * by + BLOCK_SIZE * bx;
34    C[c + wB * ty + tx] = Csub;
35 }
```