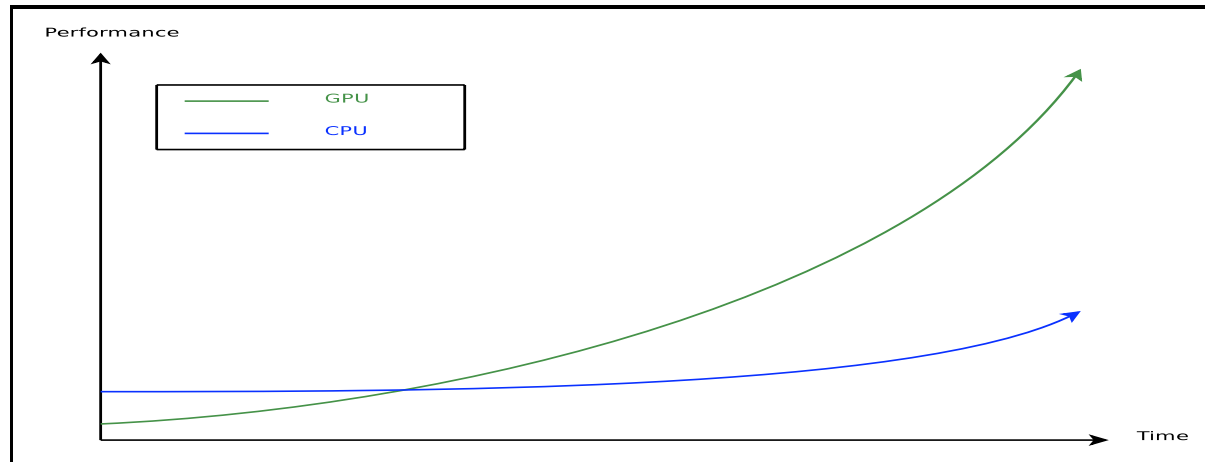

High-Performance Computing
By Advanced Stream Processing Using GPU

Se an LEE

`seanl@cse.unsw.edu.au`

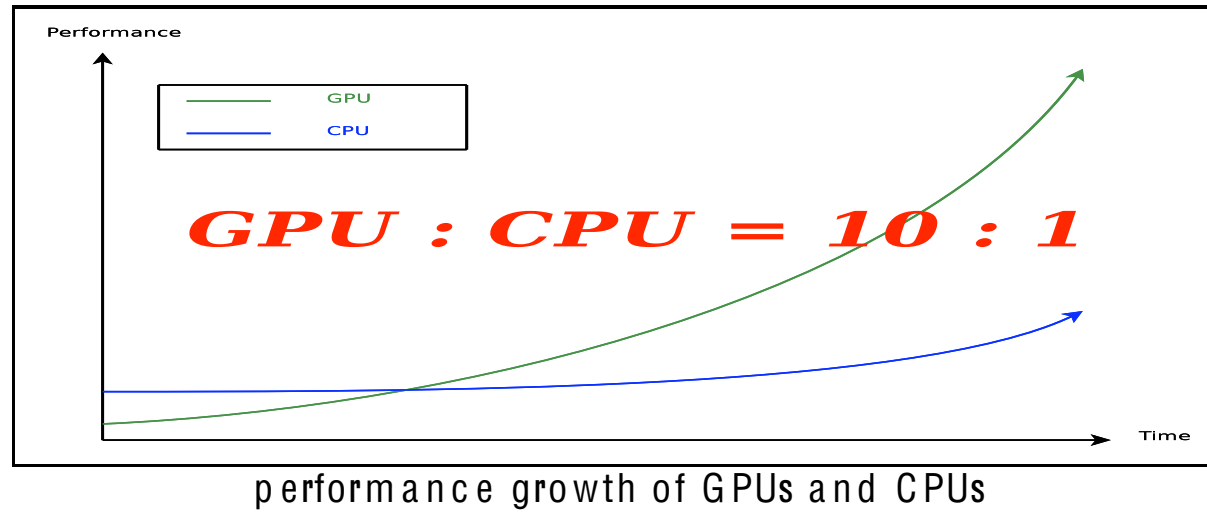
Programming Languages and Systems
School of Computer Science and Engineering
University of New South Wales

BACKGROUND

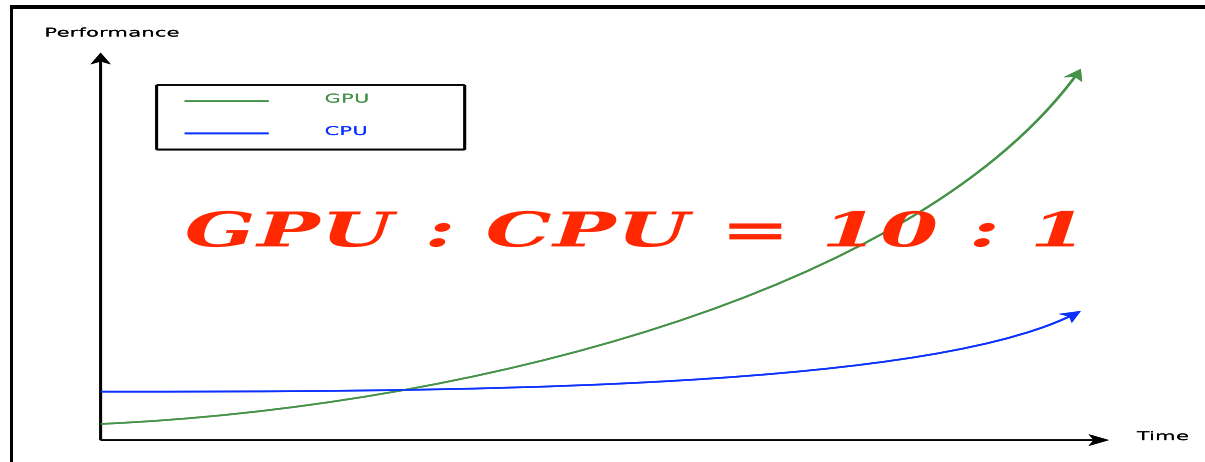


performance growth of GPUs and CPUs

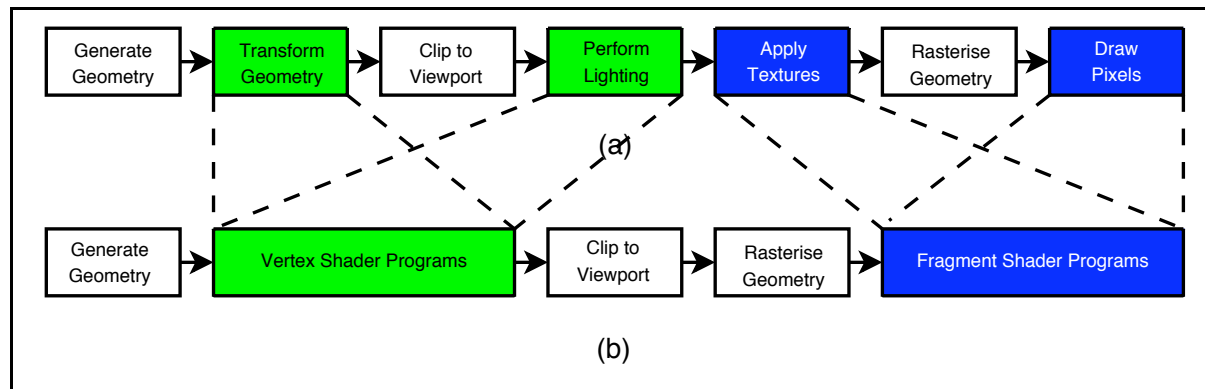
BACKGROUND



BACKGROUND



performance growth of GPUs and CPUs



(a) fixed non-programmable pipelines and (b) programmable pipelines

WHAT DO I WANT TO DO WITH GPU?

WHAT DO I WANT TO DO WITH GPU?

NDP (Nested Data Parallelism) on GPU.

WHAT DO I WANT TO DO WITH GPU?

NDP (Nested Data Parallelism) on GPU.

Why?

- GPUs cannot handle nested or irregular data structure such as sparse matrices and nested arrays. Their data structure is texture-based.
- There are algorithms that can be efficiently implemented only with nested data structure.
- Manual transformation of nested data structure is error-prone and a tedious task.

WHAT DO I WANT TO DO WITH GPU?

NDP (Nested Data Parallelism) on GPU.

In other words:

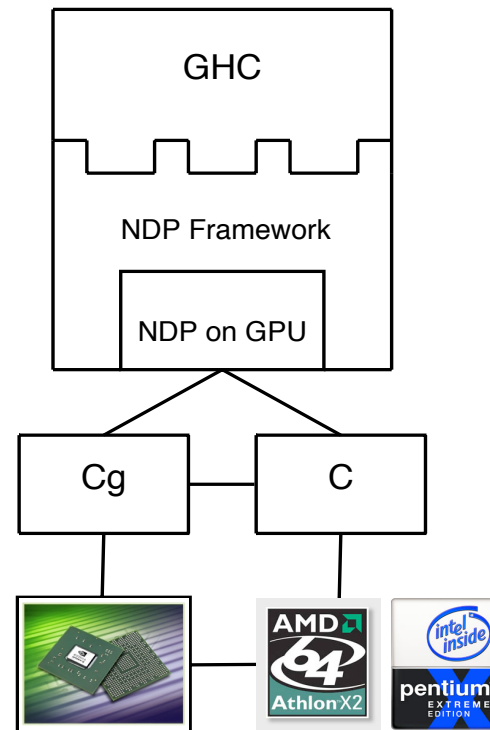
- Enhance the programmability of graphics hardware to the state where:
 - the *irregular data processing* is supported,
 - the compiler distinguishes the code to run on GPUs and the code to run on CPUs (*uniform approach*), and
 - the *hardware virtualisation* is enabled.
- Smooth integration of the enhanced programmability into the graphical applications as well as general-purpose computations.
- Achieve high-performance stream processing in low cost.

WHAT DO I WANT TO DO WITH GPU?

Other Programming Systems on GPU:

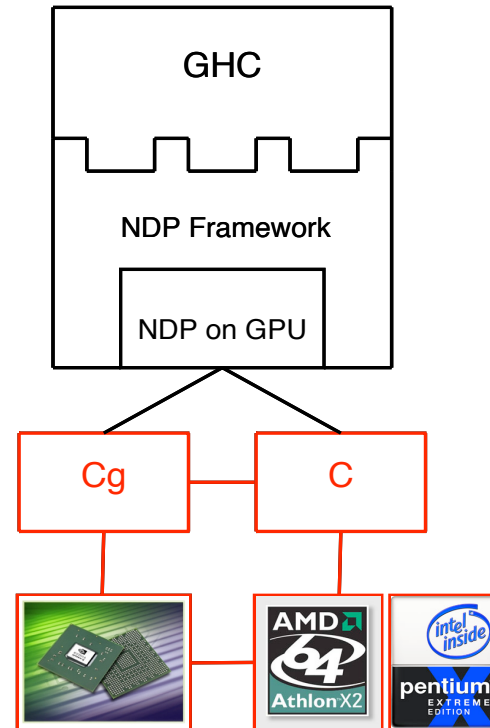
	Brook	Cg	GLSL	Sh
Nested Data Parallelism	✗	✗	✗	✗
C/C++-like syntax	✓	✓	✓	✓
Level of Abstraction	high	Low	Low	High
Standalone Compiler	✓	✓	✗	✓
Uniform Approach	✗	✗	✗	✗
Hardware Virtualisation	✓	✗	✗	✗
Irregular Data Processing Support	✗	✗	✗	✗

THE APPROACH



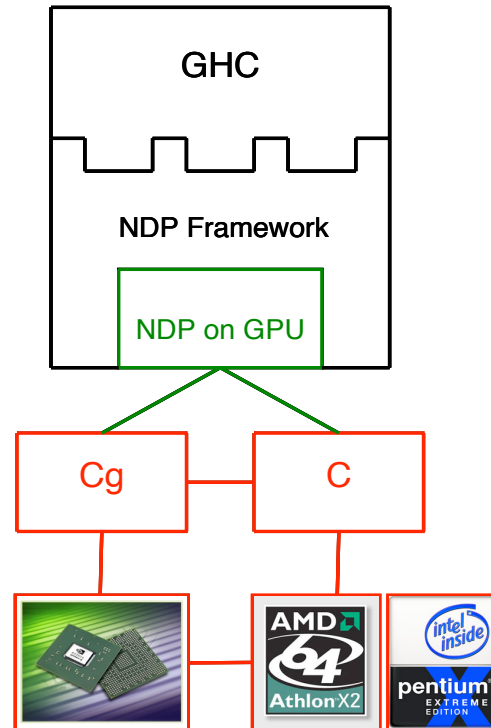
THE APPROACH

- Implement a C library for NDP on GPU. The library will include all the primitives required NDP operations on GPU.



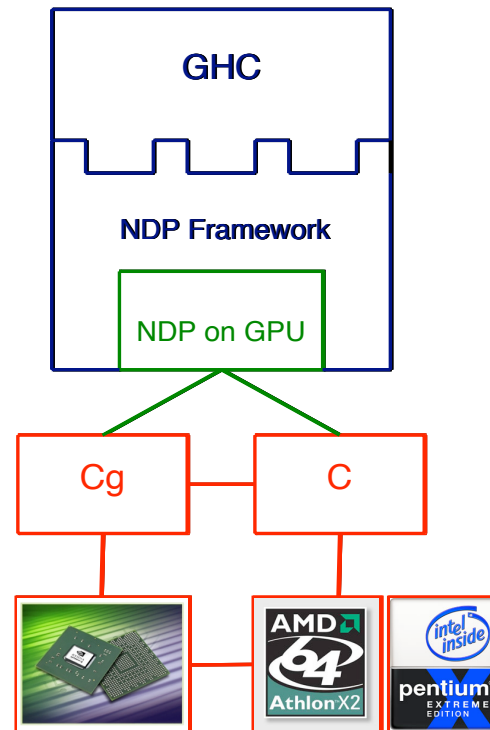
THE APPROACH

- Implement a C library for NDP on GPU. The library will include all the primitives required NDP operations on GPU.
- Develop techniques to tackle the issues such as hardware virtualisation and uniform approach.

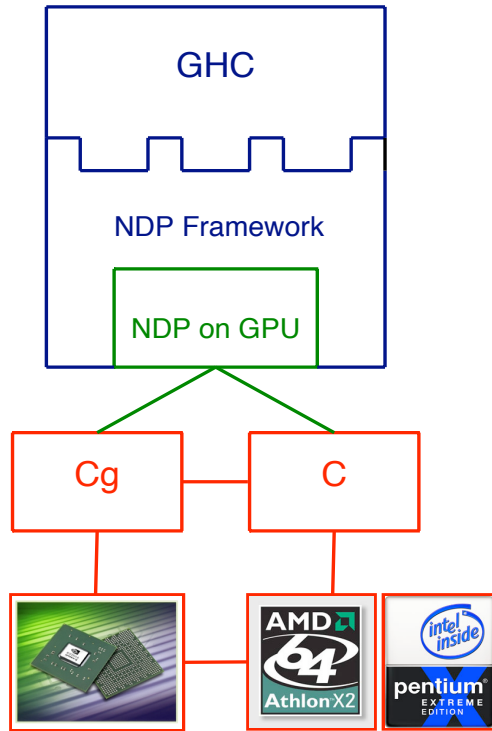


THE APPROACH

- Implement a C library for NDP on GPU. The library will include all the primitives required NDP operations on GPU.
- Develop techniques to tackle the issues such as hardware virtualisation and uniform approach.
- The library will be hooked up to NDP framework, which is being implemented in GHC.



THE APPROACH



On Haskell:

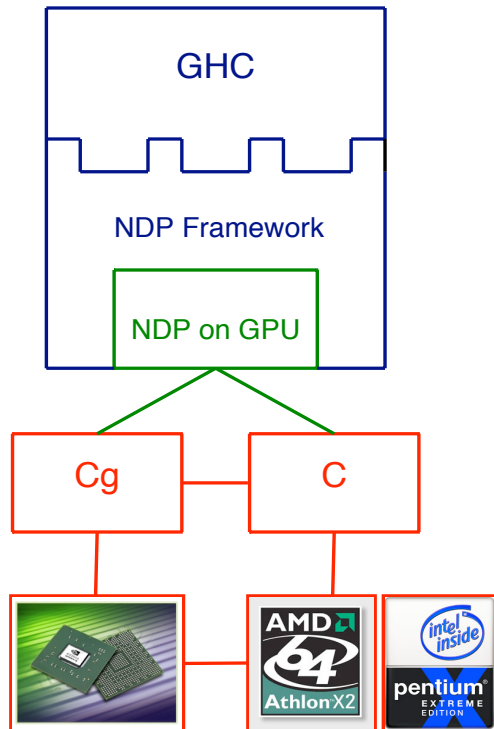
```
filter ( \x->x 'mod' 2 == 1 ) input
```

On C:

```
filter ( [ \x->x 'mod' 2 == 1 ],  
        input, inputsize, output  
);
```

THE APPROACH

On Cg:



```
float4 main ( ... data from the host CPU ... ) : COLOR {
    ... local data declaration ...
    result.a = 0;
    index[0] = ( lW * floor ( coord.y ) + floor ( coord.x ) ) * s;
    for ( i = 0; i < s ; i++, index[0]++ ) {
        if ( index[0] < numElems ) {
            index[1] = index[0];
            index[2] = 0;
            while ( index[1] >= fbWidth ) {
                index[1] -= fbWidth;
                index[2]++;
            }
            V1 = ( texRECT ( elements, float2 ( index[1], index[2] ) ) ).r;
            if ( int ( V1 ) % 2 == 1 ) {
                if ( result.a == 0 ) result.r = V1;
                else if ( result.a == 1 ) result.g = V1;
                else if ( result.a == 2 ) result.b = V1;
                result.a++;
            }
        }
    }
    return result;
}
```

WHERE AM I?

libndpgpu:

Implemented operations:

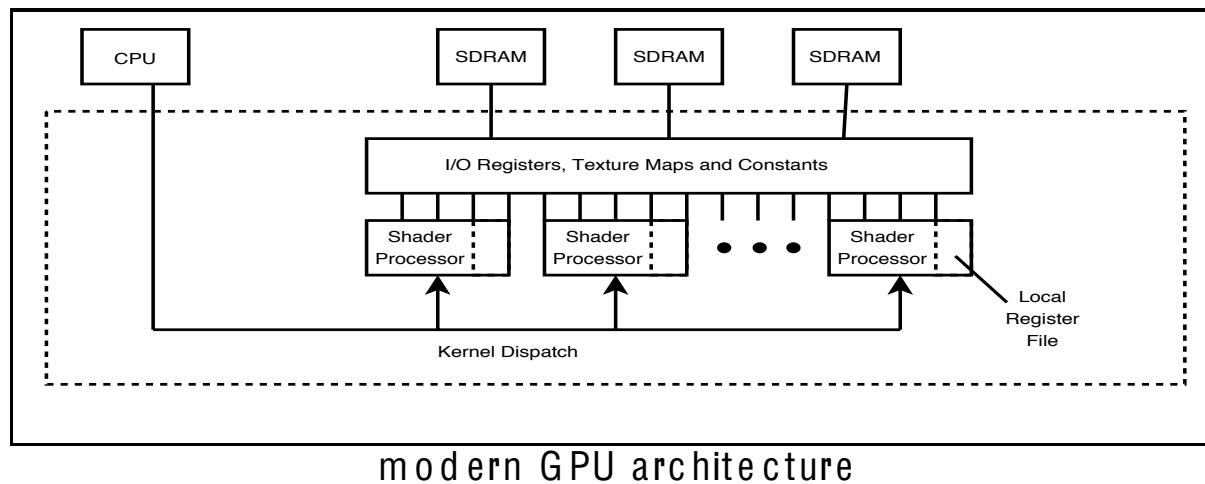
- scanl, scanlS, scanr, scanrS
- foldl, foldlS, foldr, foldrS
- map
- filter

WHERE AM I?

libndpgpu:

Implemented operations:

- scanl, scanlS, scanr, scanrS
- foldl, foldlS, foldr, foldrS
- map
- filter



WHERE AM I?

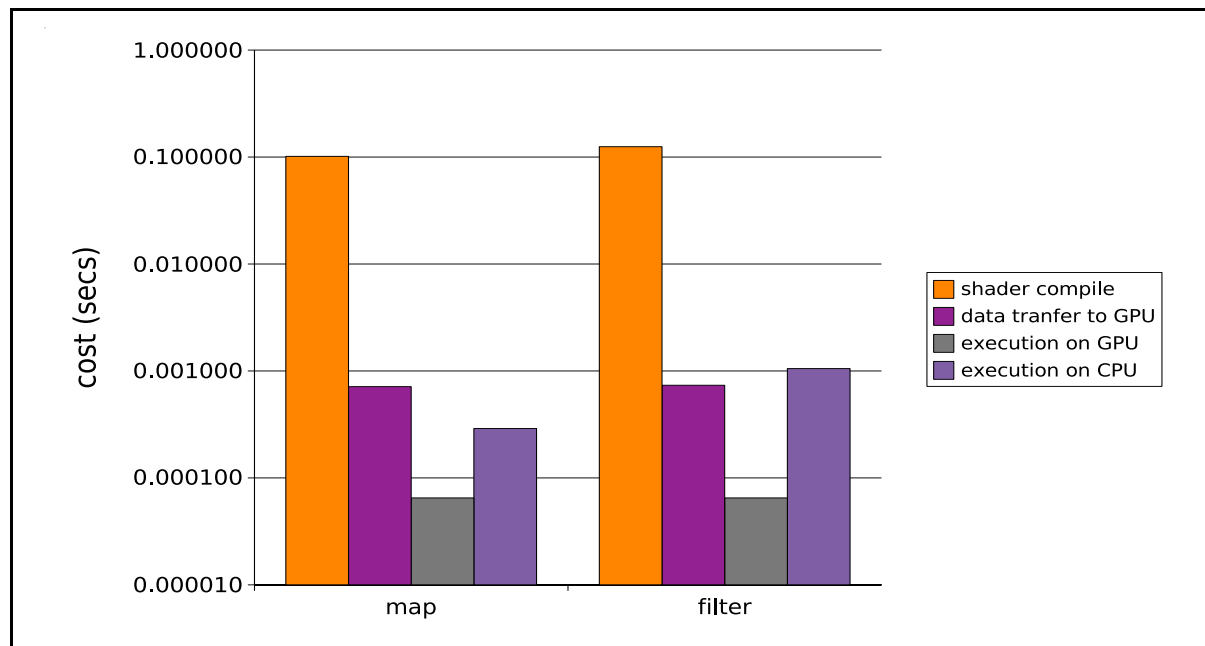
Benchmarks:

- Intel Core Duo 2.0GHz, 1 GB DDR2 RAM
- NVIDIA GeForce Go 7400 TurboCache 256MB

WHERE AM I?

Benchmarks:

- Intel Core Duo 2.0GHz, 1 GB DDR2 RAM
- NVIDIA GeForce Go 7400 TurboCache 256MB



CHALLENGES AND FUTURE WORKS (SHORT TERM)

CHALLENGES AND FUTURE WORKS (SHORT TERM)

Pre processor:

- Profile based dynamic Cg assembly generation.
- Assembly level optimisation.
- Translation of functions passed from Haskell level to lower level as inputs to higher-order functions.

CHALLENGES AND FUTURE WORKS (SHORT TERM)

Pre processor:

- Profile based dynamic Cg assembly generation.
- Assembly level optimisation.
- Translation of functions passed from Haskell level to lower level as inputs to higher-order functions.

Data organisation and transfer:

- Reduce the data transfer cost.
- Handle arrays whose number of elements are greater than the maximum size of textures using multi-pass strategy.

CHALLENGES AND FUTURE WORKS (SHORT TERM)

More bulk array operations:

- permute, zip and its variations, zipWidth and its variations

CHALLENGES AND FUTURE WORKS (SHORT TERM)

More bulk array operations:

- permute, zip and its variations, zipWidth and its variations

Release:

- Release it as a complete library by July 2007.

QUESTIONS?