

Expression Templates and OpenCL

Uwe Bawidamann,
Marco Nehmeier

Institute of Computer Science
University of Würzburg
Germany

PPAM 2011

- 1 Introduction
- 2 Expression Templates
- 3 Implementation
 - Code Generation
 - Utilize Precompiled Kernels
- 4 Experimental Results
- 5 Summary

- 1 Introduction
- 2 Expression Templates
- 3 Implementation
- 4 Experimental Results
- 5 Summary

- A change in the design of computer architecture
 - Single-core \Rightarrow multi-core
 - Circumvent physical constraints
 - Parallelize the computation

- GPGPU
- Performance
- Highly parallel
- CUDA
 - NVIDIA
 - CUDA C, CUDA Fortran or OpenCL
- OpenCL
 - Open standard
 - C like programming language
 - CPU's, GPU's, and many other

- Supercomputer
 - 3 of the 5 fastest supercomputer hosting CUDA
 - Workstation and desktop computers
 - Standard applications supporting CUDA or OpenCL
- ⇒ Write own application

- New concepts
- New Techniques
- **New pitfalls**
 - Inappropriate data transfers
 - Misaligned memory access
 - Bank conflicts
 - Inappropriate load balancing
 - Deadlocks

- Investigate techniques and concepts
 - User-friendly libraries
 - Interface in C++
 - Hiding GPGPU paradigms
- Operator overloading
 - DSL
- Expression templates
 - Code generation
 - Optimize data transfer

- 1 Introduction
- 2 Expression Templates**
- 3 Implementation
- 4 Experimental Results
- 5 Summary

■ Operator overloading

```
Vector v1, v2, v3;
...
Vector res = v1 + v2 + v3;
```

■ Pairwise evaluation problem

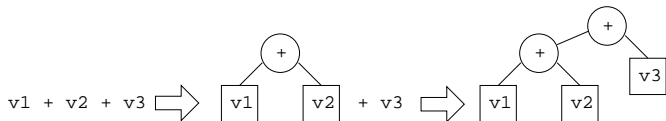
```
...
Vector tmp = v1 + v2;
Vector res = tmp + v3;
```

■ Additional loops

■ Additional instances

- new
- delete

- Introduced by Todd Veldhuizen
- Building parse trees
 - Operator overloading
 - Recursive templates
- Evaluate whole tree
 - Assignment
 - operator=
 - Optimized
 - Algorithm
 - Compiler



```
template <typename T, class OP, typename A, typename B>
class BExpr {
    A const& a_;
    B const& b_;
public:
    BExpr(A const& a, B const& b) : a_(a), b_(b) { }

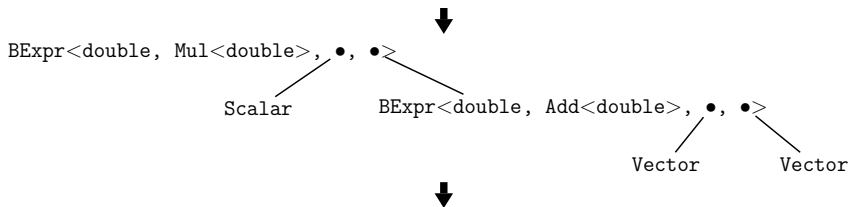
    T operator[] (size_t i) const {
        return OP.eval(a_[i], b_[i]);
    }
    ...
};
```

- Abstract representation of a binary operation
 - OP : policy class specifying operation
 - A and B : arguments
 - T : type of vector elements
- Overloaded subscript operator
 - Compute i^{th} element
 - Policy class

```
template<typename T> struct Add {  
    static T eval(T a, T b) { return a + b; }  
};
```

- Policy class
- Element-wise operation
 - Type T

Scalar * (Vector + Vector)



```

for (int i = 0; i < a.size(); ++i)
    r[i] = a[i] * (b[i] + c[i]);
  
```

- 1 Introduction
- 2 Expression Templates
- 3 Implementation**
- 4 Experimental Results
- 5 Summary

- Interval arithmetic
 - Expression templates
 - Template meta programming
- DSL
- Utilize GPU's
 - Specific Algorithms
 - Problem Solver

- Case study
 - Vector and matrix operations
- Two different approaches
 - Code generation
 - Utilize precompiled kernels

- Problems
 - Data transfer
 - Mapping
 - Code generation

- `cl::Buffer`
- Allocate and fill at creation
 - Only one data transfer
 - Permanent use of memory
- Allocate and fill within assignment operator
 - Data transfer for every expression
 - Uses only required memory

- Unique Mapping between data and kernel parameters
- Singleton pattern
 - Generate unique id's
- Objects
 - `id_`
 - `getIdent() → ''v'' + id_`
 - `getCode() → ''v'' + id_ + ''[index]''`
- Scalar
 - Value

- Computation of required parameters
 - paramSet : std::set
 - Recursively filled
 - Vectors
 - Matrices
- Header

```
std::string code = "__kernel void CLETGenVectorKernel (
for ( it = paramSet.begin() ; it != paramSet.end(); it++ ) {
    code += "__global float* " +
        (*it).getObject().getIdent() += ", ";
}
code += "const unsigned int size" +
    "\n{\n" +
    "const int index = get_global_id (0); \n";
```

■ Body

```
code += "␣" + result.getCode() + "␣=␣" +
        expression.getCode() + ";\n};\n␣\n";
```

■ getCode()

- Result type → id
- Expression tree
 - Replacement of subscript operator

```
std::string getCode() const {
    return std::string("(" + a_.getCode() +
        "␣+␣" + b_.getCode() + ")" );
}
```

- Precompiled Kernels
- Minimize data transfer
 - DSL
 - Operator overloading
- User-friendly

- Expression templates are almost similar
- Inner nodes are annotated with id's
 - Mapping between operation and temporaries
- Evaluation
 - Two traversals

- First traversal
 - Memory allocation
 - `cl::Buffer`
 - Leafs as well as inner nodes
 - `std::map`
- Second traversal
 - Computation
 - Recursive

- 1 Introduction
- 2 Expression Templates
- 3 Implementation
- 4 Experimental Results**
- 5 Summary

- Code generation vs operator overloading
 - Without compilation:
 - Speedup of about 2
 - With compilation:
 - Only for hard problems
 - Similar expressions
 - Code mixin

■ Precompiled kernels vs operator overloading

Expression	Precompiled kernels	Operator overloading
$(M1 + M2) + (M3 + M4)$	26 ms	51 ms
$(M1 + M2) + (M3 + M1)$	21 ms	51 ms
$(M1 * M2) * (M3 * M4)$	183 ms	218 ms

- 1 Introduction
- 2 Expression Templates
- 3 Implementation
- 4 Experimental Results
- 5 Summary**

- DSL
- User-friendly
- Code generation
 - Only for hard problems
 - Similar expressions
 - Code mixin
- Reduce data transfer

Questions ?

Marco Nehmeier

nehmeier@informatik.uni-wuerzburg.de