# Intel Threading Building Blocks

**Outfitting C++ for Multi-core Processor Parallelism**

Michael Müller

**dkfz.** DEUTSCHES KREBSFORSCHUNGSZENTRUM IN DER HELMHOLTZ-GEMEINSCHAFT

→Do not use threads, **automatically** map logical parallelism onto threads in a way that makes efficient use of processor resources

→Use pure C++ generic programming

→Compiler independence, processor independence, OS independence

Intel® Threading Building Blocks for Open Source

(intel)

- C++ Library for parallel computing

- Cross platform

- V4.0 since November 2011

- GPLv2 licensed, commercially available

- http://threadingbuildingblocks.org/

- TBB implements "task stealing" to balance a parallel workload

- TBB, like the STL, uses templates extensively

- TBB implements the pipeline pattern

- Basic algorithms: parallel_for, parallel_reduce, parallel_scan

- Advanced algorithms: parallel_while, parallel_do, parallel_pipeline, parallel_sort

- Containers: concurrent_queue, concurrent_vector, concurrent_hash_map

- Scalable memory allocation: scalable_malloc, scalable_free, scalable_realloc, scalable_calloc, scalable_allocator, cache_aligned_allocator

- Mutual exclusion: mutex, spin_mutex, queuing_mutex, spin_rw_mutex, queuing_rw_mutex, recursive mutex

- Atomic operations: fetch_and_add, fetch_and_increment, fetch_and_decrement, compare_and_swap, fetch_and_store

- Timing: portable fine grained global time stamp

- Task Scheduler: direct access to control the creation and activation of tasks

Serial:

```
void SerialApplyFoo( float a[], size_t n ) {
    for( size_t i=0; i<n; ++i )
        Foo(a[i]);
}
```

Parallel:

```
#include "tbb/parallel_for.h"

void ParallelApplyFoo( float a[], size_t n ) {
    parallel_for(blocked_range<size_t>(0,n,YouPickAGrainSize), ApplyFoo(a) );
}
```

Using functors:

```cpp
#include "tbb/blocked_range.h"

class ApplyFoo {
    float *const my_a;
public:
    void operator()( const blocked_range<size_t>& r ) const {
        float *a = my_a;
        for( size_t i=r.begin(); i!=r.end(); ++i )
            Foo(a[i]);
    }
    ApplyFoo( float a[] ) :
        my_a(a)
    {}
};
```

# Code Example (III)

```cpp
#include <tbb/task_scheduler_init.h>
int main() {
  tbb::task_scheduler_init init;
  // use of an algorithm
  ParallelApplyFoo( … )
}
```

Intel® Threading Building Blocks 1.0 2D Ray Tracing Application

# Using OpenMP vs. Threading Building Blocks
# for Medical Imaging on Multi-cores

Philipp Kegel, Maraike Schellmann, and Sergei Gorlatch

University of Münster, Germany
{p.kegel,schellmann,gorlatch}@uni-muenster.de

**Abstract.** We compare two parallel programming approaches for multi-core systems: the well-known OpenMP and the recently introduced Threading Building Blocks (TBB) library by Intel®. The comparison is made using the parallelization of a real-world numerical algorithm for medical imaging. We develop several parallel implementations, and compare them w.r.t. programming effort, programming style and abstraction, and runtime performance. We show that TBB requires a considerable program re-design, whereas with OpenMP simple compiler directives are sufficient. While TBB appears to be less appropriate for parallelizing existing implementations, it fosters a good programming style and higher abstraction level for newly developed parallel programs. Our experimental measurements on a dual quad-core system demonstrate that OpenMP slightly outperforms TBB in our implementation.

Michael Müller
E130

4/10/2012 | Seite 11

**TBB vs. OpenMP: What the official FAQ says**

**dkfz.**

→ **What about OpenMP?**

Everyone should use OpenMP as much as they can. It is easy to use, it is standard, it is supported by all major compilers, and it exploits parallelism well. But it is very loop oriented, and does not address algorithm or data structure level parallelism.

- **TBB**: C++ library for parallel computing: Focuses on tasks not threads, nice programming style

- **OpenMP**: Same basic concept, but not C++, loop oriented

- **ITK, Qt Threads**: Abstraction layers to Windows Threads or Unix pthreads