# Leveraging Non-Blocking Collective Communication in High-Performance Applications

Torsten Hoefler, Peter Gottschling and Andrew Lumsdaine

Open Systems Lab
Indiana University
Bloomington, USA

20th ACM Symposium on Parallelism in Algorithms and Architectures - SPAA'08

Munich, Germany

June, 14th 2008

## Features of non-blocking collective operations

- hide full communication latency by overlapping
- use the available bandwidth better
- avoid detrimental effects of pseudo-synchronization/process skew
- make efficient use of the new semantics

## LibNBC and MPI

- implements all MPI collectives non-blocking
- overhead-optimized implementation
- special InfiniBand$^{TM}$ optimizations
- progress thread

# Problems and a Solution

## Challenges for the Programer

- rearrange the algorithm to overlap
- implement and debug non-blocking communication
- optimize overlap (e.g., message sizes)

## Overcoming the Problems

- semi-automatic approach for applications with independent data
- covers many applications that fit the map-reduce model
- many scientific applications (e.g., parallel data processing, Fourier transformation, parallel sorting, FEM methods, ...)

# A typical Program - Parallel Compression

```
1   my_size = 0;
2   for (i=0; i < N/P; i++) {
3     my_size += compress(i, outptr);
4     outptr += my_size;
5   }
6   gather(sizes, my_size);
7   gatherv(outbuf, sizes);
```

```
1  for (i=0; i < N/P; i++) {
2    my_size = compress(i, outptr);
3    gather(sizes, my_size);
4    igatherv(outptr, sizes, hndl[i]);
5    outptr += my_size;
6    if(i>0) waitall(hndl[i-1], 1);
7  }
8  waitall(hndl[N/P], 1);
```

```
1    for (i=0; i < N/P/t; i++) {
2      size = 0;
3      for (j=i; j < i+t; j++) {
4        my_size = compress(i*t+j, outptr);
5        outptr += my_size;
6        size += my_size;
7      }
8      gather(sizes, size);
9      igatherv(outptr-size, sizes, hndl[i]);
10     if(i>0) waitall(hndl[i-1], 1);
11   }
12   waitall(hndl[N/P/t], 1);
```

# Parallel Compression - Adding a Window

```
1   for ( i =0; i < N/P/ t ; i ++) {
2     my_size = 0;
3     for ( j=i ; j < i+t ; j ++) {
4       my_size += compress ( i ∗t+j , outptr );
5       outptr += my_size ;
6     }
7     gather ( sizes , my_size );
8     igather ( outbuf , sizes , hndl [ i ] );
9     if ( i > w) waitall ( hndl [ i −w] , 1);
10  }
11  waitall ( hnld [N/P/ t −w] , w);
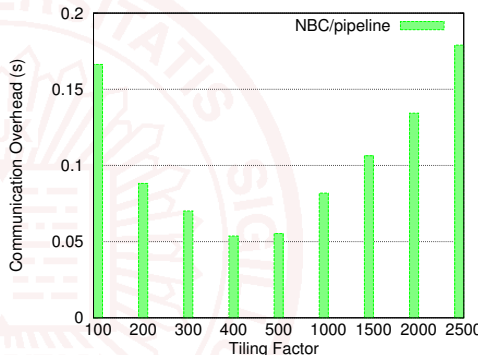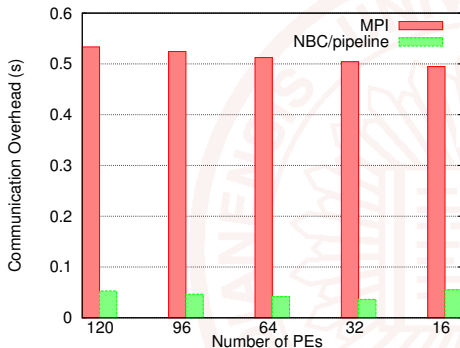```

# Automatic Transformation

## Templated Transformation

- requires buffer, computation and communication functor
- C++ template tiles loops and uses window
- $\Rightarrow$ programmer-directed overlap simplifies optimization
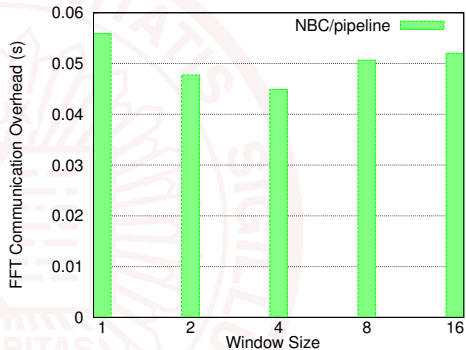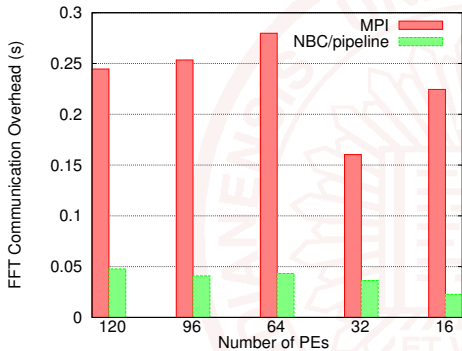


## Two Examples

- parallel compression
- parallel 3d Fast Fourier Transformation

- 128 2 GHz Opteron 246 nodes, InfiniBand[TM]
- 146MiB data compressed with `bzip2`
- 21% speedup on 120 PEs

- 16% speedup on 120 PEs
- weak scaling ($400^3, 480^3, \ldots, 720^3$

## Conclusions

- loop-tiling and introduction of a commmunication-window to leverage non-blocking operations
- proposed a template-driven optimization scheme to assist the programmer
- showed the usefulness and performance advantages with two applications
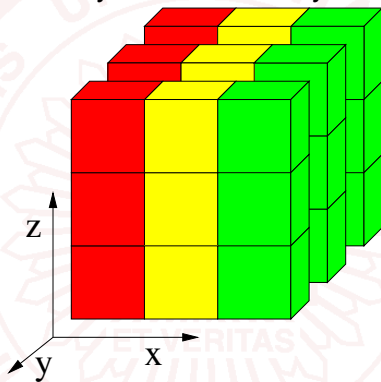- LibNBC and templates available at:
  `http://www.unixer.de/NBC`

## Future Work

- optimize more (real-world) applications
- automatic parameter tuning

# Conclusions and Future Work

## Conclusions

- loop-tiling and introduction of a commmunication-window to leverage non-blocking operations
- proposed a template-driven optimization scheme to assist the programmer
- showed the usefulness and performance advantages with two applications
- LibNBC and templates available at:
  `http://www.unixer.de/NBC`

## Future Work

- optimize more (real-world) applications
- automatic parameter tuning

# Backup Slides

Data already transformed in y direction

1 block = 1 double value (3x3x3 grid)
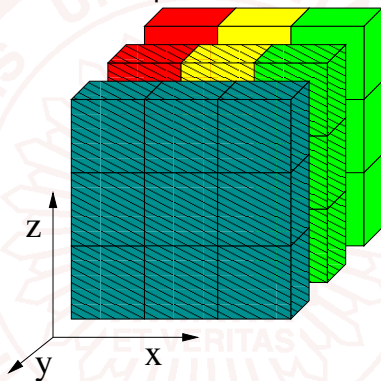
Transform first xz plane in z direction



pattern means that data was transformed in y and z direction
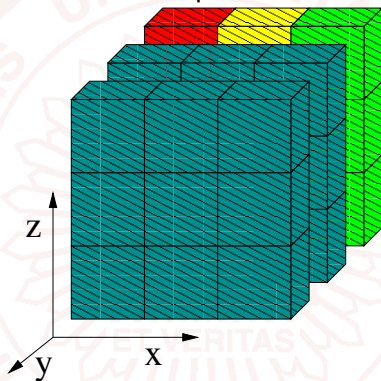
start MPI_Ialltoall of first xz plane and transform second plane
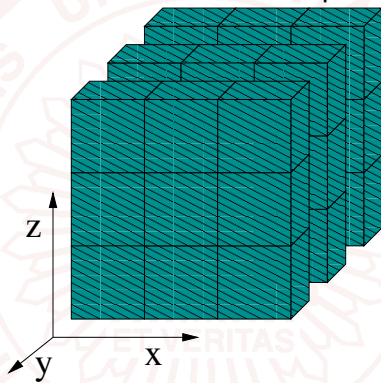


cyan color means that data is communicated in the background

start MPI_Ialltoall of second xz plane and transform third plane



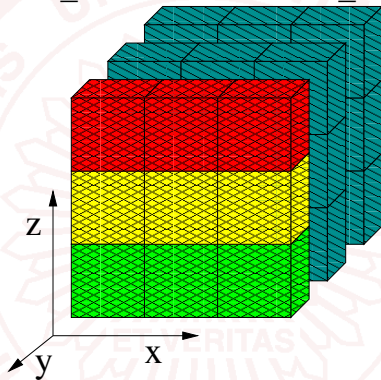data of two planes is not accessible due to communication

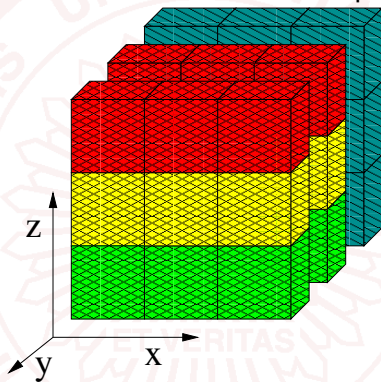start communication of the third plane and ...



we need the first xz plane to go on ...

... so MPI_Wait for the first MPI_Ialltoall!



and transform first plane (new pattern means xyz transformed)

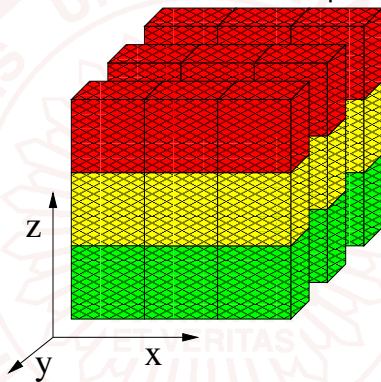Wait and transform second xz plane



first plane's data could be accessed for next operation

wait and transform last xz plane



done! $\rightarrow$ 1 complete 1D-FFT overlaps a communication